

Machine Learning for Engineering Problem Solving

A Practical Example-driven Guide to
Classical Techniques



Austin R.J. Downey

Department of Mechanical Engineering
Department of Civil and Environmental Engineering
University of South Carolina, Columbia SC, USA

July 21, 2025

Contents

Preface	1
Accompanying Video Lectures	1
Programming in Python	1
Cover Art	2
License	2
Source Code	2
Questions and Contact Information	2
1 Basic Concepts in Machine Learning	3
1.1 Examples of Artificial Intelligence	4
1.1.1 Examples of Artificial Intelligence	4
1.1.2 Examples of Machine Learning	5
1.1.3 Examples of Deep Learning	6
1.2 Definition of Machine Learning	9
1.3 Supervision in Machine Learning	10
1.3.1 Supervised Learning	10
1.3.2 Unsupervised Learning	11
1.3.3 Semisupervised Learning	13
1.3.4 Reinforcement Learning	13
1.4 Types of learning (Batch and Online)	14
1.4.1 Batch Learning	14
1.4.2 Online Learning	14
1.4.3 Bad Data in Machine Learning	15
1.5 Learning Approaches: Instance-Based and Model-Based	16
1.5.1 Instance-based learning	16
1.5.2 Model-based learning	17
1.5.3 Selecting methods:	18
1.6 The Unreasonable Effectiveness of Data	18
2 Regression	19
2.1 Linear Regression	21
2.1.1 Closed Form Solution	22
2.1.2 Computational Complexity	24
2.1.3 Comparison of Gradient Descent Methods	27
2.1.4 Batch Gradient Descent	27
2.1.5 Stochastic Gradient Descent	30
2.1.6 Mini-batch Gradient Descent	31
2.1.7 Feature Scaling	31
2.2 Polynomial Regression	33
2.3 Training and Testing Data	35
2.4 Pipelines	36
2.5 Learning Curves	36
2.6 Regularized Linear Models	39

2.6.1	Ridge Regression	39
2.7	Early Stopping	41
2.8	Examples	42
3	Classification	52
3.1	Binary Classifier	52
3.1.1	Regularized Linear Classifier	54
3.2	Performance Measures for Binary Classification	55
3.2.1	Confusion Matrix	55
3.2.2	Accuracy, Precision, and Recall	57
3.3	k -fold Cross-validation	61
3.4	Multiclass Classification	62
3.5	Performance Measures for Multiclass Classification	64
3.5.1	Confusion Matrix	64
3.5.2	Analyzing Individual Errors	67
3.6	Examples	68
4	Regression-Based Classification	79
4.1	Logistic Regression	79
4.1.1	1-D Decision Boundaries	82
4.1.2	2-D Decision Boundaries	83
4.2	Softmax Regression	84
4.3	Examples	87
5	Decision Trees	94
5.1	Decision Tree Classification	94
5.1.1	Class Probability Estimation	96
5.2	The CART Training Algorithm	97
5.2.1	Computational Complexity	99
5.2.2	Entropy verse Gini Impurity	99
5.3	Decision Tree Regression	99
5.4	Instability	101
5.5	Random Forest	103
5.6	Examples	104
6	Support Vector Machines	107
6.1	Linear SVM Classification	108
6.1.1	Decision Function and Predictions	110
6.1.2	Training Objective	111
6.1.3	Quadratic Programming	112
6.2	Nonlinear SVM Classification	113
6.2.1	Kernel trick	114
6.2.2	Standard Kernels	115
6.3	Computational Complexity	116
6.4	SVM Regression	117

6.4.1	Linear SVR	117
6.4.2	Kernel SVR	118
6.5	Examples	119

Preface

This text is a practical, example-driven guide to introduce classical machine learning techniques using the `scikit-learn` library designed for engineers with limited to no programming experience. This preface collects the essential housekeeping information for using this text.

Accompanying Video Lectures

Videos of lectures associated with this text are available as a playlist [here](#).

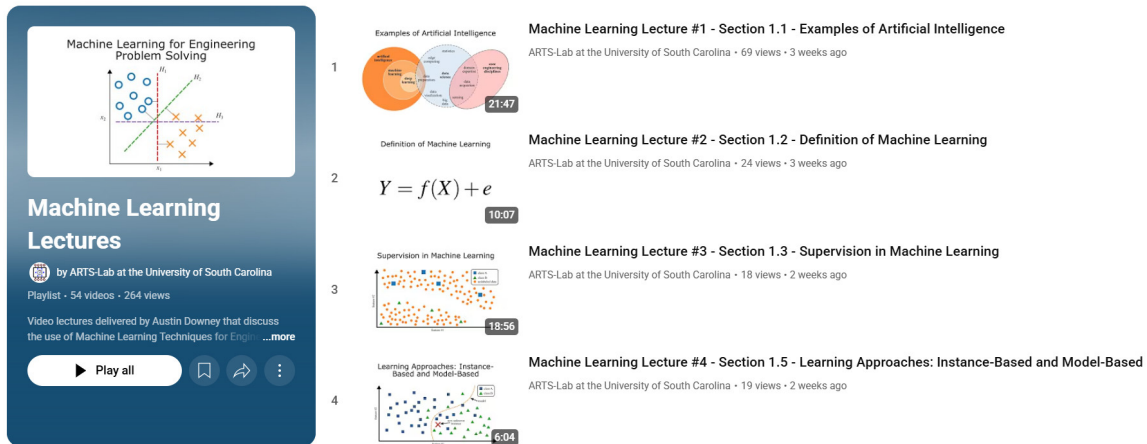


Figure 1: Playlist of videos associated with this text.

Programming in Python

This text uses Python programmed through the Spyder IDE managed through the Anaconda platform for the examples, leveraging the `scikit-learn` library to explain the topics discussed. To assist readers of the text, a six part video series that walks the practitioner through this combination of IDE and distribution manager is provided as a playlist [here](#).

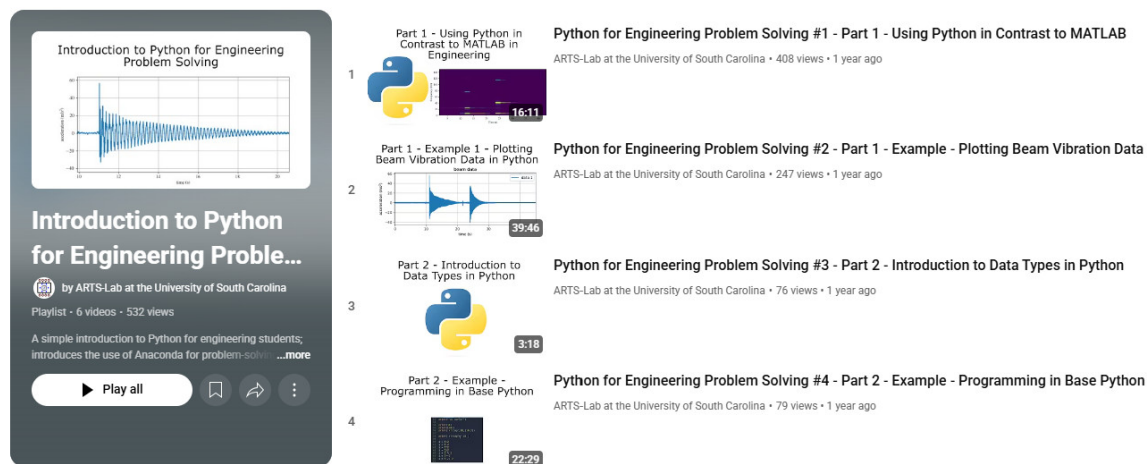


Figure 2: Playlist of videos for learning how to program in Python using Spyder and Anaconda.

Cover Art

The cover image, is a mid-1930s International Harvester C-series truck. Built between 1934 and 1936 at the company's Springfield Works in Springfield, Ohio. Roughly 80,000 C-series trucks were made during this short time. The C-series was International's first line to feature an all-steel cab and a host of mechanical upgrades over the prior W-series. This truck was photographed by Ethan McNeese on Washington Island in Door County, Wisconsin in 2021.

The truck on the cover underscores the simple idea that a machine's purpose is to turn raw input into useful work. The same principle drives machine learning, where digital "machines" transform data into insight and data-informed actions.

License

This work is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License (CC BY-SA 4.0). More information on the Attribution-ShareAlike 4.0 International license can be found [here](#). Unless otherwise denoted, all text, figures, diagrams, and photos used in this work are the sole property of the authors and are released under CC BY-SA 4.0 both in part and in whole. Reworks and redistributions of this work that fall within the CC BY-SA 4.0 licenses are encouraged.

Source Code

The source code for this text is available [here](#).

Questions and Contact Information

For questions, corrections, or requests, contact Austin Downey at austindowney@sc.edu.

1 Basic Concepts in Machine Learning

Machine Learning (ML) is a method of data analysis that automates analytical model building. ML is closely coupled to data science, as shown in figure 1.1, which is a multidisciplinary field that utilizes scientific methods, processes, algorithms, and systems to extract knowledge and insights from structured and unstructured data. It is based on the idea that systems can learn from data, identify patterns, and make decisions with minimal human intervention. This section introduces some foundational concepts in machine learning and its applications.

- **ML is not creating robots:** A common misconception is that machine learning is about building robots. In reality, ML focuses on developing algorithms that can learn from and make predictions or decisions based on data.
- **The SPAM filter is one of the initial ML uses:** One of the earliest and most common applications of ML is the spam filter, which classifies emails as spam or not spam. This has been followed by numerous other applications such as:
 - **Speech to text technology:** Converting spoken language into written text, which is used in virtual assistants and transcription services.
 - **Medical diagnostics:** Assisting doctors by predicting diseases from medical images and patient data.
- **ML has lots of fundamental concepts (jargon):** To effectively understand and apply machine learning, it's essential to grasp several key concepts and terminologies, including:
 - **Supervised vs unsupervised learning:** Supervised learning involves training a model on labeled data, while unsupervised learning deals with data that has no labels and tries to find hidden patterns.

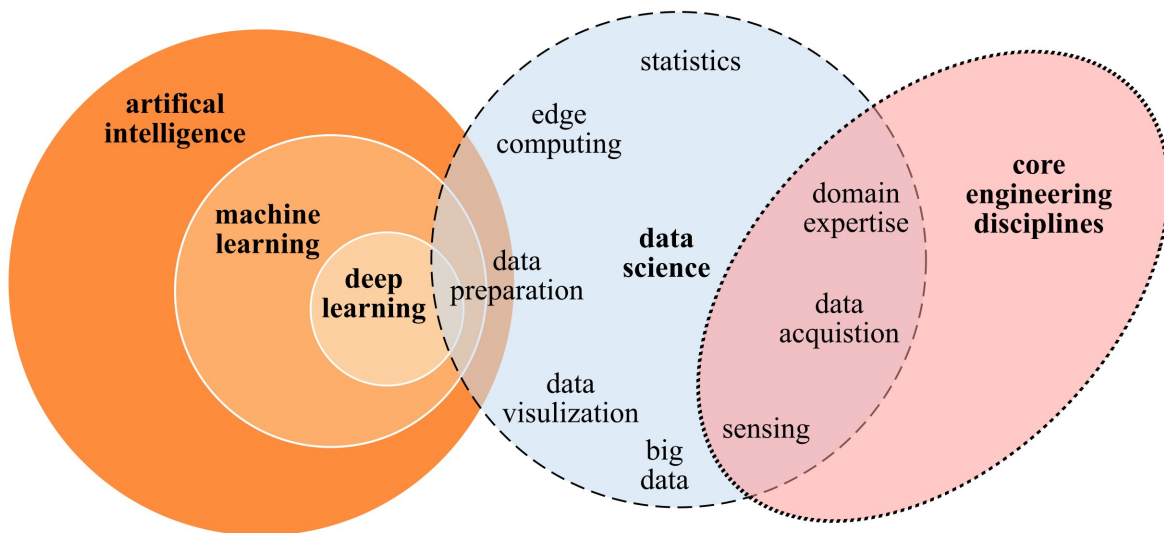


Figure 1.1: Overlap between the fields of Artificial Intelligence (AI), Data science (DS), and the core engineering disciplines of Civil, Mechanical, Electrical, and Chemical Engineering.

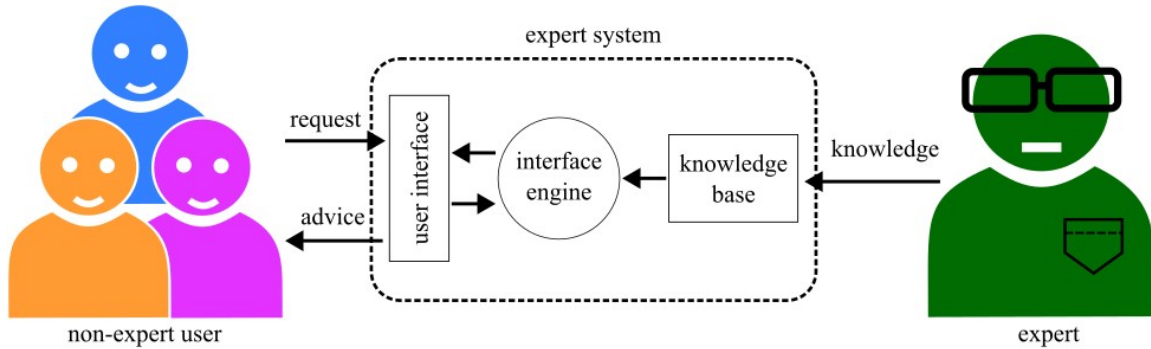


Figure 1.2: Diagram of an “expert system” in AI, which is a computer program that simulates the decision-making ability of a human expert by using a knowledge base and inference rules.

- **Online versus batch learning:** Online learning algorithms update the model incrementally as new data arrives, whereas batch learning algorithms train the model using the entire dataset at once.
- **Instance-based vs model-based learning:** Instance-based learning algorithms, such as k-nearest neighbors, use specific instances to make predictions, whereas model-based algorithms, like linear regression, build a model from the training data and use it to make predictions.

1.1 Examples of Artificial Intelligence

The field of Artificial Intelligence (AI) encompasses a diverse array of technologies and methodologies aimed at enabling machines to perform tasks that typically require human intelligence. This section explores various examples of AI, Machine Learning, and Deep Learning technologies, highlighting their distinct characteristics and applications.

1.1.1 Examples of Artificial Intelligence

Artificial Intelligence (AI) encompasses a wide range of technologies aimed at making machines simulate human intelligence. Some examples include:

- **Expert systems:** Computer programs that simulate the decision-making ability of a human expert by using a knowledge base and inference rules, as diagrammed in figure 1.2.
- **Chatbots:** Programs designed to simulate conversation with human users, especially over the internet.



Figure 1.3: A Symbolics Lisp Machine, a specialized hardware platform designed to run expert systems which are a version of AI focusing on answering questions to challenging problems.^a

1.1.2 Examples of Machine Learning

Machine Learning, as a subset of AI, involves algorithms that improve automatically through experience. Some common examples include:

- **Linear regression:** A statistical method for modeling the relationship between a dependent variable and one or more independent variables.
- **Classification:** Techniques such as decision trees and support vector machines (SVMs) that categorize data into predefined classes.
- **Simple image/speech recognition:** Algorithms that identify objects in images or convert spoken language into text, fundamental in applications like facial recognition and virtual assistants. An example is the Symbolics Lisp Machine shown in figure 1.3.

^aMichael L. Umbricht and Carl R. Friend (Retro-Computing Society of RI), CC BY-SA 3.0 <<https://creativecommons.org/licenses/by-sa/3.0/>>, via Wikimedia Commons

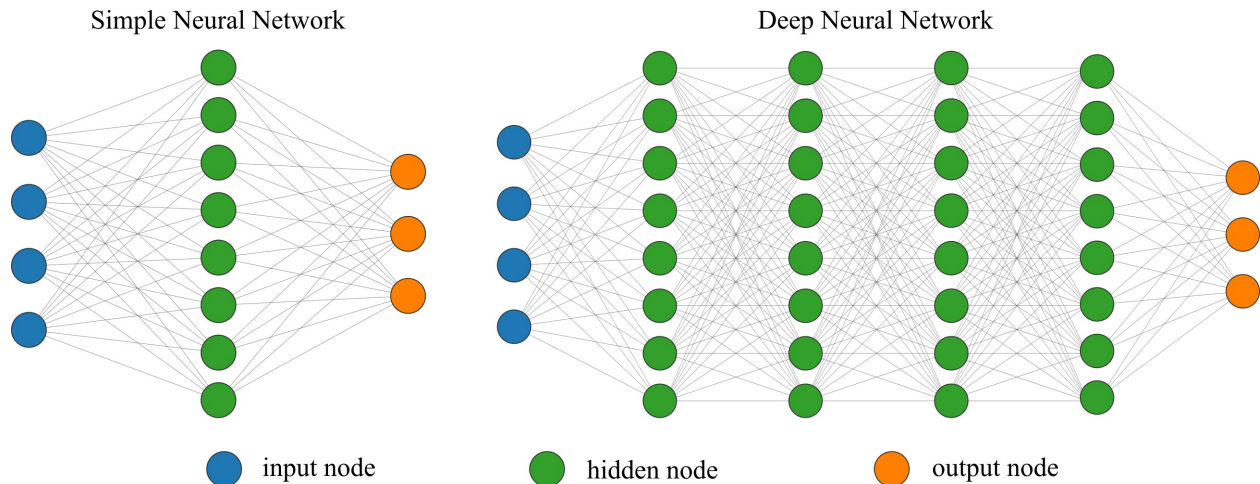


Figure 1.4: simple neural network vs deep learning

1.1.3 Examples of Deep Learning

Deep learning, an enhancement of machine learning, utilizes “deep” neural networks to construct knowledge graphs, as shown in figure 1.4. Though conceptualized in the 1970s, it was initially unfeasible due to the problem of vanishing gradients within the networks. In 2012, Geoffrey E. Hinton’s team demonstrated that a network with 60 million parameters and 650,000 neurons could effectively perform image classification across a dataset containing 1,000 categories (paper shown in figure 1.5). This breakthrough was facilitated by the use of GPUs and a novel regularization technique known as “dropout.” The team’s modified model participated in the ILSVRC-2012 competition, securing a first-place top-5 test error rate of 15.3%, a significant improvement over the 26.2% recorded by the runner-up and a major leap when compared to the previous year as diagrammed in figure 1.6.

Review 1.1 Imagenet in 2012 represented a significant step forward in machine learning by introducing the first practical example of deep learning, famously known as AlexNet (Figure 1.5) This model, developed by Geoffrey Hinton and his team, utilized deep convolutional neural networks to dramatically improve the accuracy of image classification, which was a longstanding challenge in the field.

ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca

Abstract

We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet ILSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operation. To reduce overfitting in the fully-connected layers we employed a recently-developed regularization method called “dropout” that proved to be very effective. We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.

Figure 1.5: Geoffrey’s 2012 paper “ImageNet Classification with Deep Convolutional Neural Networks”.^a

The results of AlexNet were rather monumental, reducing the top-5 test error rate to 15.3% compared to 26.2% by the next best entry, as shown in Figure 1.6. This was clear evidence of deep learning’s superior capability over traditional machine learning methods, effectively revolutionizing the approach towards machine learning in the broader scientific community.

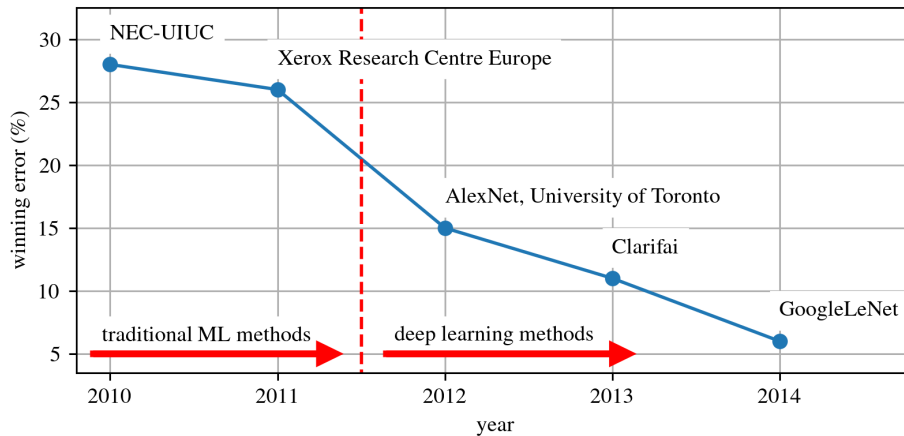


Figure 1.6: ImageNet Competition Results showing the impact of deep learning methods on image classification.

The implications of this leap forward led to broad applications of deep learning that permeate numerous aspects of technology and science today. The impact of AlexNet and subsequent deep learning developments culminated in the awarding of the 2024 Nobel Prize in Physics to Geoffrey Hinton, recognizing his contributions to the field of artificial neural networks. Figure 1.7 shows the 2024 Nobel Award Ceremony in Stockholm Sweden with John Hopfield and Geoffrey Hinton receiving their awards from the King of Sweden.

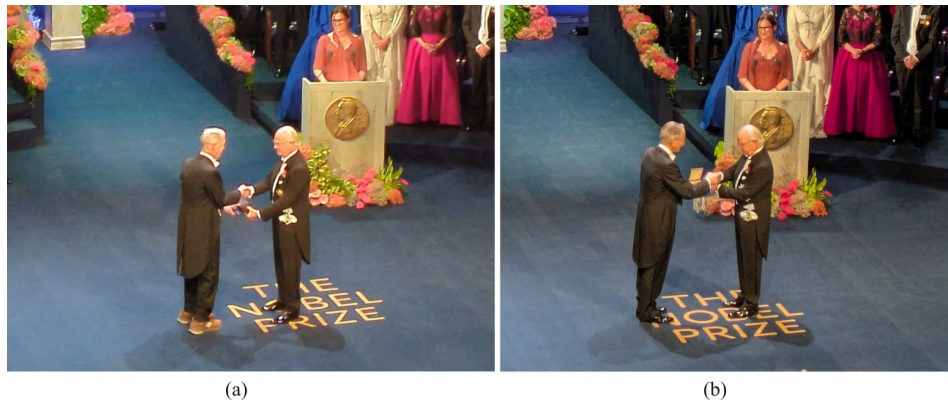


Figure 1.7: The 2024 Nobel Prize in Physics for “for foundational discoveries and inventions that enable machine learning with artificial neural networks” was awarded to: (a) John Hopfield and (b) Geoffrey Hinton. ^a

^aCopyright held by Authors and Neural Information Processing Systems Foundation, Inc., used under fair use. <<https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf>>

^aThe author of this text was lucky enough to attend the 2024 Nobel Prize Ceremony in Stockholm and took these pictures.

1.2 Definition of Machine Learning

The definitions and conventions that follow provide a common language and point of reference for all subsequent material:

- Machine learning is the discipline of creating computer programs that improve automatically by analyzing data. Here a broader and a more engineering-focused definition of Machine learning is provided⁵:
 - [Machine Learning is the] field of study that gives computers the ability to learn without being explicitly programmed. Arthur Samuel, 1959.
 - A computer program is said to learn from experience E with respect to some task T and some performance measure P , if its performance on T , as measured by P , improves with experience E . Tom Mitchell, 1997.
- The use of X and Y for variables.
 - For the i^{th} sample, $x^{(i)}$ contains all its input features (excluding the target), while $y^{(i)}$ denotes the corresponding target value.
 - General definition: “Machine learning algorithms are described as learning a target function (f) that best maps input variables (X) to an output variable (Y).”

$$Y = f(X) \tag{1.1}$$

This is described as a standard learning challenge where the goal is to predict future values Y using new samples of input variables X . The function f that relates inputs to outputs is not known. If it were, direct application would be possible, eliminating the necessity for learning it via machine learning methods. This process is more complex than it may initially seem. Furthermore, there is an error e associated with this task that is independent of the input data X .

$$Y = f(X) + e \tag{1.2}$$

This yields two primary phases in the machine-learning workflow:

- Training: Creating the model is a compute-intensive process often run in a data center
- Inference: Using the model can be computationally cheap and even performed “at the edge”

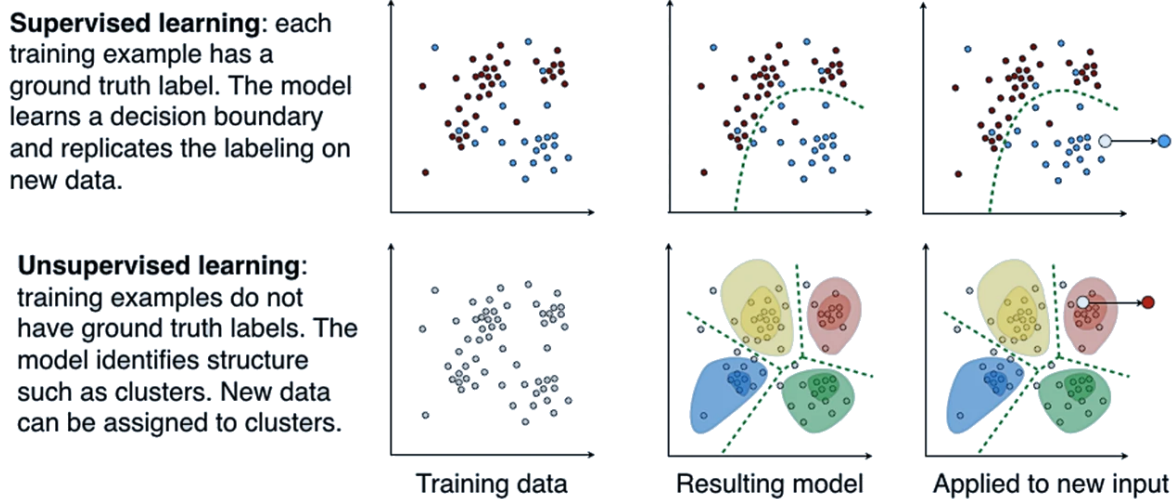


Figure 1.8: Supervised vs unsupervised machine learning methods^a.

1.3 Supervision in Machine Learning

The general domains of ML can be classified into supervised and unsupervised learning, as shown in figure 1.8.

1.3.1 Supervised Learning

In supervised learning, the data is pre-categorized with labels

- **Classification**, The process of categorizing a given set of data into classes
- **Regression**, The process of estimating the relationships between a dependent variables (i.e. output) and one or more independent variables (i.e. input).

In supervised learning, the training data provided to the algorithm includes the desired solutions, known as labels. A common task within this type of learning is classification. An example of classification is a spam filter, which is trained using numerous emails that are each labeled as either spam or ham. The objective for the spam filter is to learn how to accurately classify new emails based on this training. Regression is another typical task is to predict a target numeric value, such as the price of a house, given a set of features (size, proximity to railroad or highways, age of house, style, etc.) called predictors. To train a system you need to build a model and train in on many house sale examples, including both predictors (house features) and their labels (i.e., prices). Some of the most important supervised learning algorithms:

- k-Nearest Neighbors
- Linear Regression
- Logistic Regression
- Support Vector Machines (SVMs)
- Decision Trees and Random Forests

^amodified from: Langs, G., Röhrich, S., Hofmanninger, J. et al., CC BY 4.0 <<https://creativecommons.org/licenses/by/4.0/>>, via Wikimedia Commons

- **Clustering**, The process of identifying and grouping similar data points in larger datasets without concern for the specific outcome
- **Association**, The process learning a rule-based method for discovering relations between variables data data
- **Dimension Reduction**, The process of reducing the number of input variables in training data.

Here are some of the most important unsupervised learning algorithms:

- Clustering
- k-Means
- Hierarchical Cluster Analysis (HCA)
- Expectation Maximization

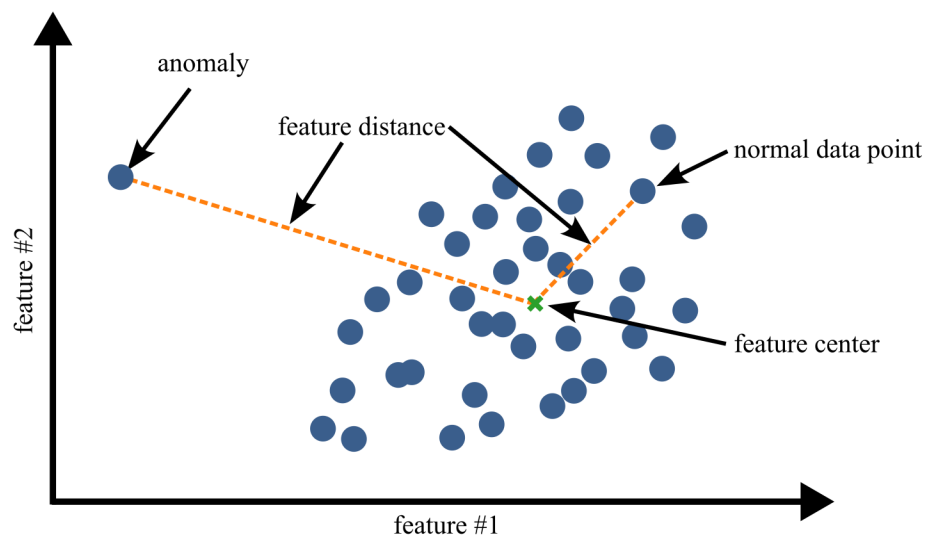


Figure 1.10: Anomaly Detection for a given set of data.

Another key application of unsupervised learning is anomaly detection as shown in figure 1.10. The goal is to flag unusual credit-card transactions that may signal fraud, spot defects on a production line, or filter out outliers before passing a dataset to another learning algorithm. The model is first exposed to many examples of normal behavior so it can build a reference profile. When a new observation arrives it checks how closely the example matches that profile and labels it normal or anomalous accordingly.

1.3.3 Semisupervised Learning

Certain algorithms are designed to manage partially labeled training data, typically characterized by a substantial amount of unlabeled data with a minimal amount of labeled data as shown in figure 1.11. Many semi-supervised learning algorithms are hybrid forms, integrating features of both unsupervised and supervised learning approaches.

An example is seen in image hosting services like Google Photos. Upon uploading your family images, the system automatically recognizes and groups repeated appearances of the same individuals across different photos, such as person 1 in photos A, C, and D and person 2 in photos A, B, and E. This clustering is the unsupervised component of the algorithm. The system then only needs a single label for each person to subsequently identify and tag everyone in all images, enhancing the ease of searching through photos.

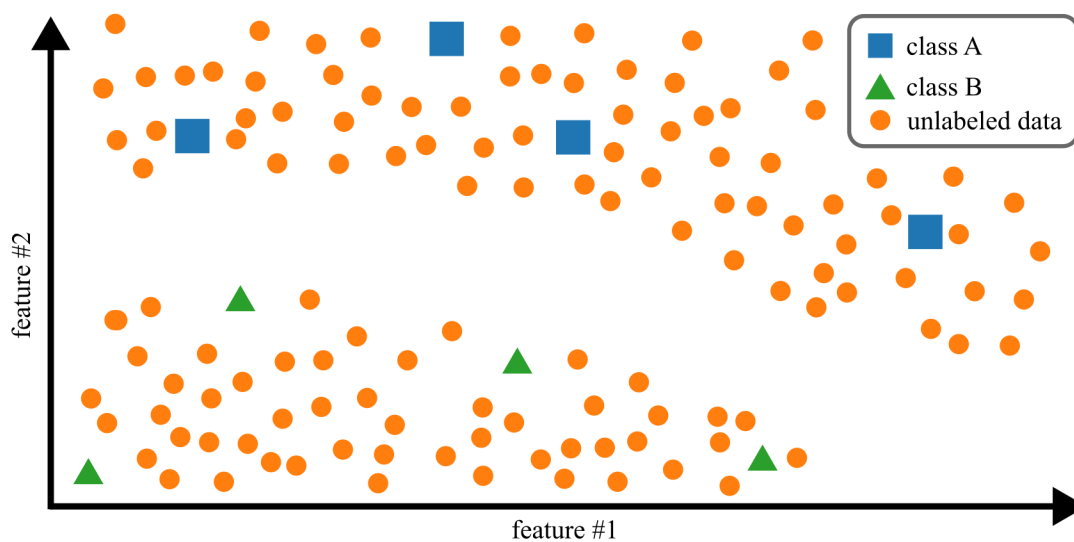


Figure 1.11: Semisupervised learning enabling the use of a limited set of labeled data to infer the labels of larger unlabeled datasets.

1.3.4 Reinforcement Learning

Reinforcement learning follows a distinct paradigm in which an autonomous agent repeatedly interacts with its environment. At each step the agent observes the current state, chooses and executes an action, and then receives feedback in the form of a reward or penalty; as shown in figure 1.12. By sampling many such cycles the agent gradually discovers a strategy, known as a policy, that maps perceived situations to actions so as to maximise the total reward accumulated over time.

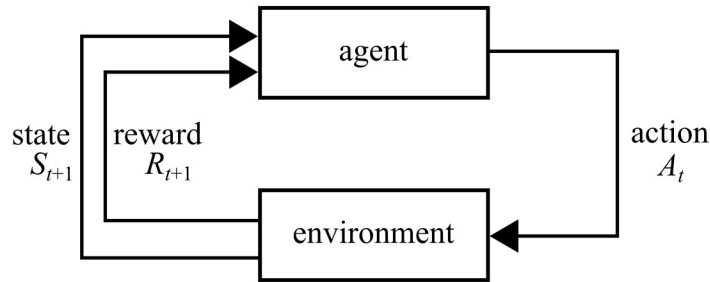


Figure 1.12: reinforcement learning control diagram

1.4 Types of learning (Batch and Online)

Machine learning systems can be broadly categorized based on how they process and learn from data.

1.4.1 Batch Learning

In batch learning, the system is fully trained using the entirety of the available data, which often demands significant time and computational resources, and is thus usually performed offline. Initially, the system undergoes training; after which, it is deployed into production where it no longer learns but merely applies the previously acquired knowledge. This method is referred to as offline learning. The deployment phase, known as inference, is typically executed swiftly.

To update a batch learning system with new information, such as recognizing a novel form of spam, it is necessary to develop a completely new version of the system. This version must be trained from the ground up using the entire dataset, both old and new data, before replacing the older system in production.

Despite these challenges, the training, evaluation, and deployment processes of a machine learning system can be automated, allowing even batch learning systems to adapt to changes. This approach is straightforward and usually effective; however, training on a full dataset can be time-consuming - often taking many hours - hence, systems are usually updated no more frequently than daily or weekly. Moreover, utilizing the full dataset requires extensive computing resources, such as CPU power, memory, disk space, and network bandwidth. For organizations with vast amounts of data, the costs of daily retraining from scratch can be prohibitively expensive.

If the dataset is exceptionally large, employing a batch learning algorithm may become impractical. Additionally, in situations where autonomous learning is essential and computational resources are limited, such as with smartphone applications or extraterrestrial rovers, the need to manage large datasets and conduct lengthy training sessions daily poses significant challenges.

1.4.2 Online Learning

Online learning trains a model incrementally by presenting new data points one at a time or in small mini-batches. Each update is computationally light and fast, allowing the model to refresh its knowledge continuously as fresh data streams in. This approach is particularly appropriate for systems that:

- receive data continuously, such as stock prices,
- need to quickly or autonomously adapt to changes,
- have limited computing resources: once an online learning system processes new data instances, they can be discarded to save space, unless there is a need to revert to a previous state and “replay” the data.

Online learning is also useful for managing large datasets that exceed the memory capacity of a single machine, known as out-of-core learning. The algorithm processes parts of the data, conducts a training step, and repeats this until all the data has been processed.

A critical parameter in online learning systems is the learning rate, which dictates how rapidly the system adapts to changing data. A high learning rate allows for rapid adaptation but may also lead to quick forgetting of old data and training on noise. A low learning rate results in slower learning and reduced sensitivity to variations in new data.

1.4.3 Bad Data in Machine Learning

A significant challenge in online learning is the system’s susceptibility to performance degradation when exposed to poor quality data. This is particularly problematic in live systems where clients might quickly notice issues. For instance, bad data could originate from a malfunctioning robot sensor or an attempt to manipulate search engine rankings through spamming. To mitigate these risks, it is crucial to vigilantly monitor the system and quickly disable learning or revert to a previously effective state if a decline in performance is observed. Additionally, monitoring the input data for anomalies and employing anomaly detection algorithms can help identify and respond to aberrant data.

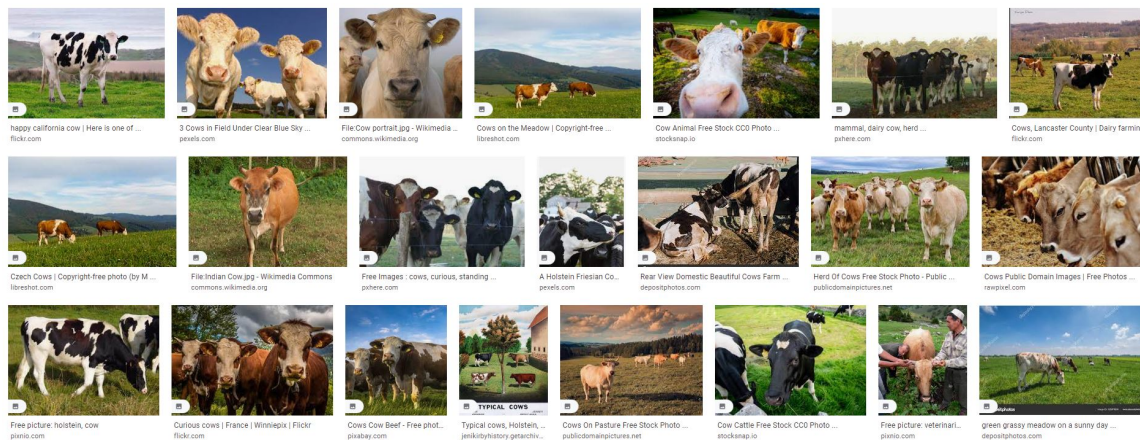


Figure 1.13: Training a machine learning algorithm to recognize cows may end up just learning to recognize grass. This is sometimes called “short-cut learning”^a.

For example, imagine training a convolutional neural network to recognize cows (Figure 1.13). Every image in the training set shows black-and-white cattle standing on lush green pasture, so the

^aScreen shoot of a Google image search for cows taken by the Austin Downey, all images stated to be under a creative commons license per Google search tools; also assumed to be fair use under given the educational purpose of this text, via <https://www.google.com/>

easiest statistical cue for the model to latch onto is the dominant green background rather than the animals themselves. At inference time, the network confidently labels any scene filled with grass as “cow,” yet fails when presented with a cow on snow or asphalt. In other words, it has learned “grass recognition,” not “cow recognition”.

1.5 Learning Approaches: Instance-Based and Model-Based

Another method of classifying machine learning systems is based on their generalization capabilities. Typically, the main objective in machine learning is to make predictions, which requires the system to generalize from its training examples to new, unseen instances. Achieving high performance on training data is useful but not the ultimate goal; the system must also excel when confronted with new data. There are primarily two strategies for generalization: instance-based learning and model-based learning.

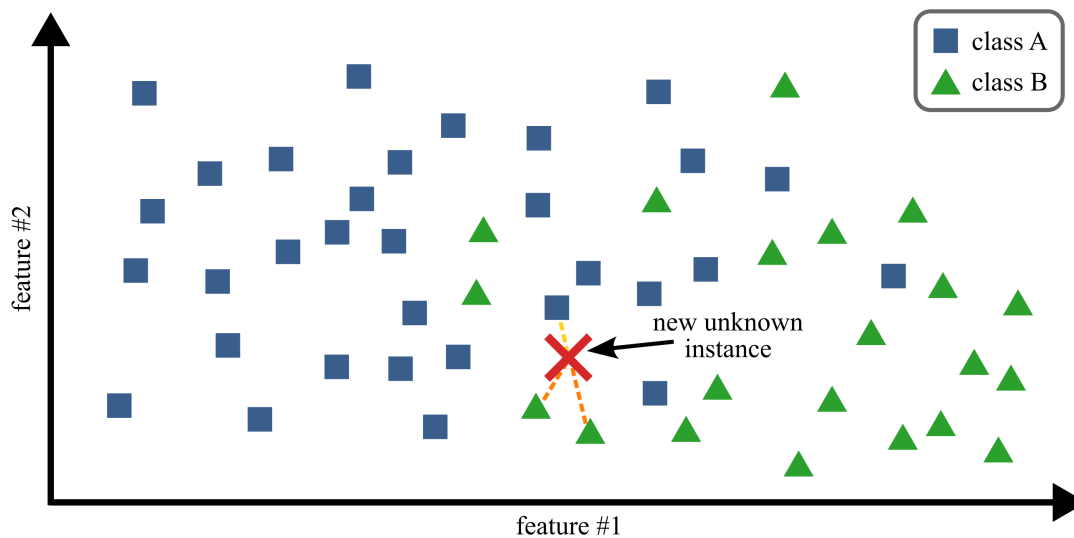


Figure 1.14: Instance-based learning where the class of a unknown instance is inferred from the its distance to data points with known labels.

1.5.1 Instance-based learning

One of the simplest learning strategies is rote memorization. For instance, a spam filter built on this principle would only mark emails as spam if they match previously flagged emails exactly. While this approach is straightforward, it’s hardly the most effective. A more sophisticated spam filter could extend its detection capabilities to include emails that closely resemble known spam messages. This approach necessitates a metric for measuring the similarity between two emails. A basic method might involve counting the shared words between emails. Under this system, an email would be classified as spam if it shares a significant number of words with an email already identified as spam. This method exemplifies instance-based learning (Figure 1.14), where the system memorizes specific examples and uses a similarity metric to generalize to new instances.

1.5.2 Model-based learning

Another approach to generalization involves constructing a model based on a set of examples and then using this model to make predictions as shown in figure 1.15. This method is known as model-based learning.

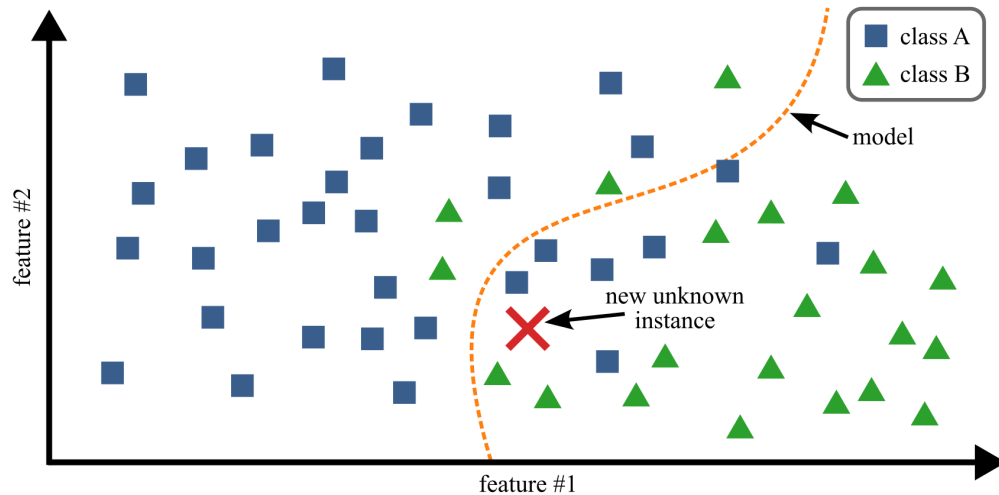


Figure 1.15: Model-based learning where the class of a unknown instance is inferred from its location in reference to a model trained on the data with known labels.

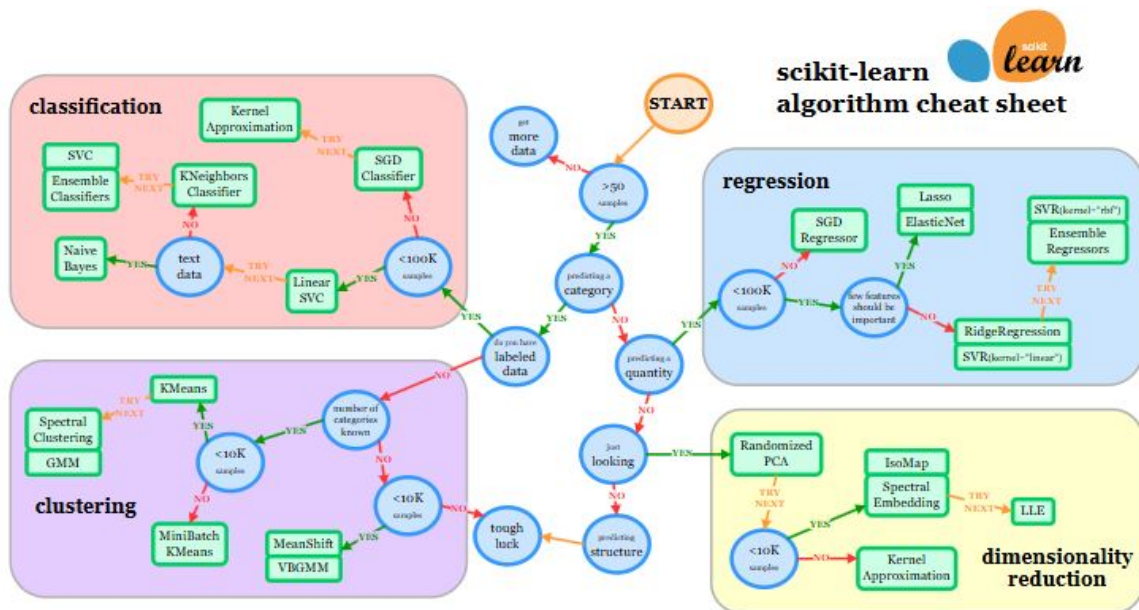


Figure 1.16: Flowchart of estimators used in the scikit learn library that intends guide users for what algorithms to use for a given case^a.

^aScikit-learn algorithms cheat sheet, permissive simplified BSD license and assumed to be fair use under given its nature as documentation and the educational purpose of this text, via https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html

1.5.3 Selecting methods:

One of the initial challenges faced by newcomers to machine learning is selecting the appropriate algorithms (or estimators) for specific tasks. This is due to the fact that various algorithms excel with different kinds of data and problems. The flowchart depicted in Fig. 1.16, which originates from the scikit-learn library documentation, is intended to guide users in choosing the most suitable algorithms for their particular datasets.

1.6 The Unreasonable Effectiveness of Data

In a landmark study released in 2001, Microsoft researchers Michele Banko and Eric Brill demonstrated that a variety of Machine Learning algorithms, even relatively straightforward ones, achieved similar levels of performance on a complex task of natural language disambiguation when provided with sufficient data. This finding is illustrated in figure 1.17 and further discussed in their paper^a.

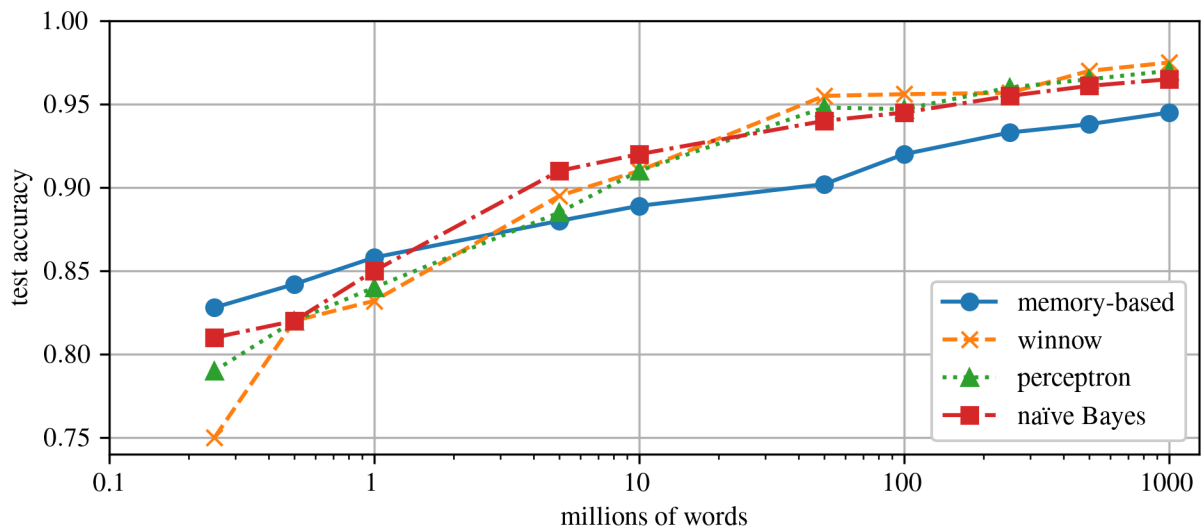


Figure 1.17: Learning curves for four algorithms studied in Banko and Brill showing the importance of the amount of data when compared to algorithm selection.

As the authors put it: “these results suggest that we may want to reconsider the tradeoff between spending time and money on algorithm development versus spending it on corpus development.”. Or put more directly “The Unreasonable Effectiveness of Data”^b. It’s important to recognize, however, that small and medium-sized datasets remain prevalent, and acquiring additional training data is not always straightforward or economical. Therefore, the significance of algorithm selection should not be overlooked.

^aMichele Banko and Eric Brill. “Scaling to very very large corpora for natural language disambiguation.” Proceedings of the 39th annual meeting of the Association for Computational Linguistics. 2001.

^bHalevy, Alon, Peter Norvig, and Fernando Pereira. “The unreasonable effectiveness of data.” IEEE intelligent systems 24.2 (2009): 8-12.

2 Regression

Regression is a fundamental tool in machine learning used to model the relationship between an independent variable (called data and denoted x) and a dependent variable (called target and denoted y). The goal is to learn a function that best predicts the target value from the input data. This relationship is shown in figure 2.1 .

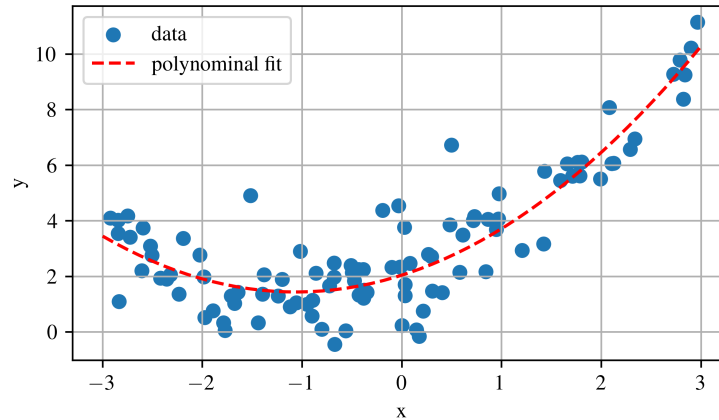


Figure 2.1: Polynomial model fit via regression to a generic dataset.

In this chapter, we first examine Linear Regression and contrast two ways to estimate its parameters:

1. **Closed-form solution.** Solve for the parameter vector in one step by minimizing the cost function on the entire training set.
2. **Gradient Descent .** Apply an iterative optimizer that repeatedly updates the parameters in small steps that lower the cost until it reaches the same optimum found by the closed-form method. We will look at three common Gradient Descent variants: Batch Gradient Descent, Mini-batch Gradient Descent, and Stochastic Gradient Descent.

The first approach gives an exact answer immediately, whereas the second reaches that answer through successive refinements.

Review 2.1 Ames Housing Dataset

The Ames housing dataset provides details on individual residential property sales in Ames, Iowa, from 2006 to 2010 ^a. It includes 2,930 entries and numerous variables (23 nominal, 23 ordinal, 14 discrete, and 20 continuous) that are used to evaluate home values ^b.

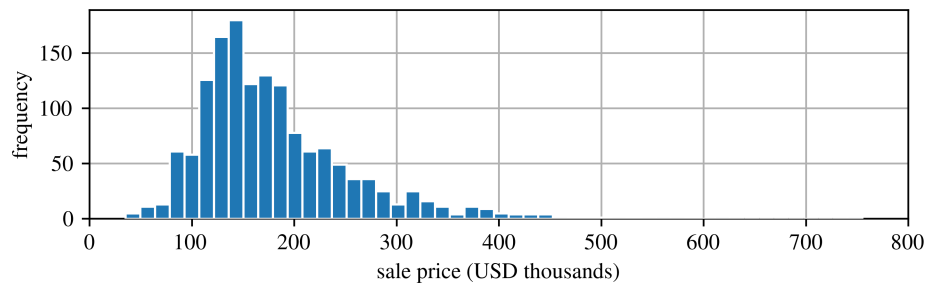


Figure 2.2: Histogram of sale prices from 2006 to 2010 for the 2,930 entries.

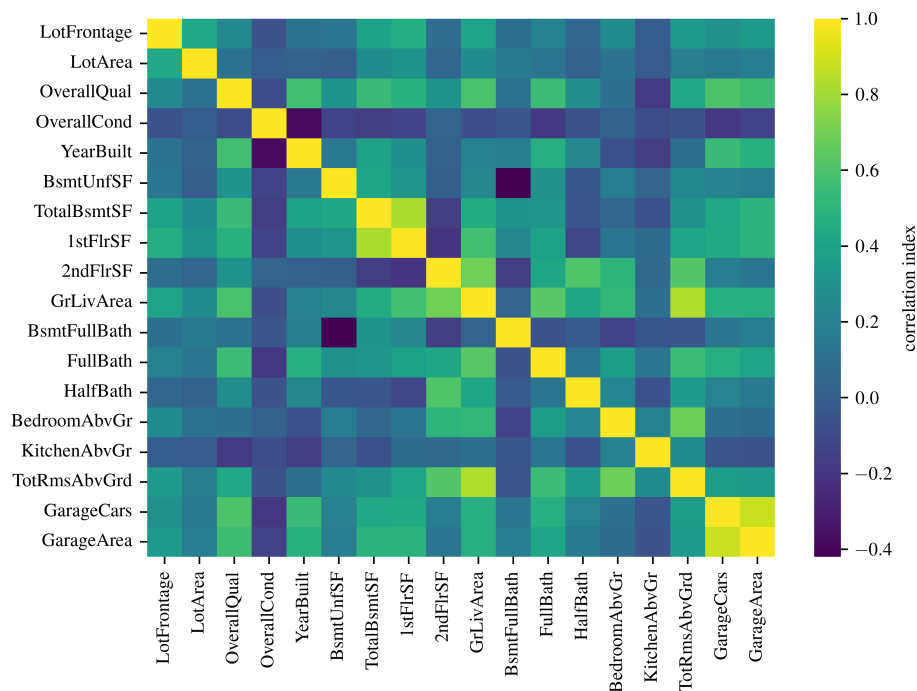


Figure 2.3: Correlation matrix for the 13 input parameters.

^aThe author of this text bought a home in Ames Iowa during this time right as the 2008 housing crash unfolded and on the last day they allowed home loans with no down payment.

^bDe Cock, Dean. "Ames, Iowa: Alternative to the Boston housing data as an end of semester regression project." Journal of Statistics Education 19.3 (2011).

2.1 Linear Regression

Consider the scenario where you wish to determine the value of a house in Ames using the Ames dataset included in sklearn. To illustrate this, let us plot the relationship between the above ground living area and the housing price.

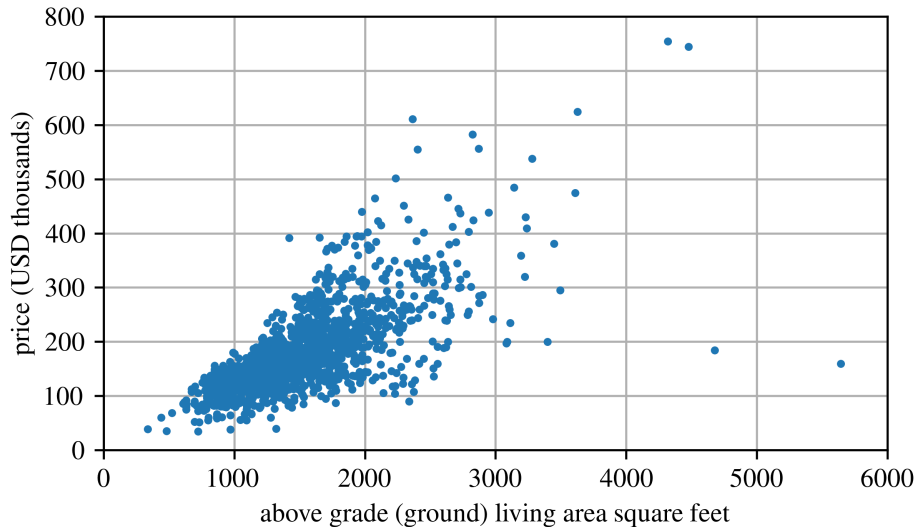


Figure 2.4: Ames housing data.

A clear trend is observable, despite the data's noise (i.e., partial randomness): the value of a house increases with the above ground living area. Therefore, the house value can be modeled as a linear function of the above ground living area:

$$\text{price_model} = \theta_1 + \theta_2 \times \text{above_ground_living_area} \quad (2.1)$$

This model comprises two parameters, θ_1 and θ_2 , which can be adjusted to represent any linear relationship.

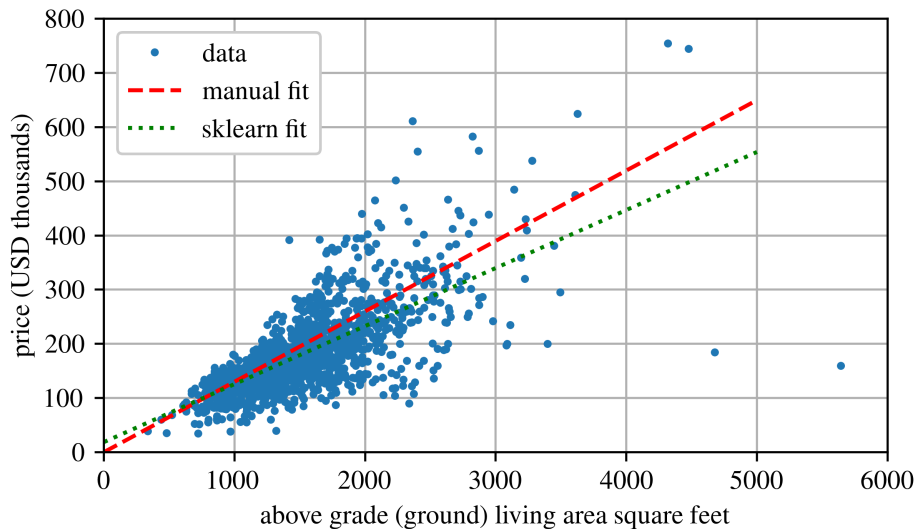


Figure 2.5: Ames housing data with trend lines.

Notations

- \mathbf{x} : a lowercase, bolded variable that denotes a vector.
- \mathbf{X} : an uppercase, bolded variable that represents a matrix.
- \mathbf{X} : this denotes the feature matrix, also known as the input matrix.
- \mathbf{y} : this is the label vector corresponding to the feature matrix, representing the target outputs.
- $\mathbf{x}^{(i)}$ represents the i^{th} instance in the dataset \mathbf{X} . For instance, a dataset about cats might include an entry for a specific cat named Mittens. This entry (labelled y) could detail several attributes such as age, weight, and length.
- h : the hypothesis or prediction function that, given a feature vector \mathbf{X} , outputs the predicted values $\hat{\mathbf{y}} = h(\mathbf{X})$.
- n : the total number of features in the dataset, which defines the complexity or degree of the model.
- m : the total number of instances in the dataset.
- \hat{x} : an estimated value, indicated by a hat, and pronounced as “ x -hat”.

2.1.1 Closed Form Solution

Before you can tune a model’s parameters, you need a clear way to judge how well it fits the data. The standard approach is to define a cost function that assigns a numerical penalty to poor fits; lower values indicate better performance. For a linear model, this cost function measures the difference between the model’s predictions and the actual training targets, and you minimize that difference using Ordinary Least Squares linear regression. Practically, you feed your training data into the algorithm, which then solves for the parameter values that align the linear relationship most closely to the data. In Python, this procedure is implemented in `scikit-learn` as `sklearn.linear_model.LinearRegression`.

In the context of the Ames housing data, a straightforward cost function used is the Mean Squared Error (MSE), defined as

$$\text{MSE} = J = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \quad (2.2)$$

where m is the number of dataset instances, and J serves as a general representation of the cost function in machine learning contexts.

The MSE is also applied to minimize the error in a linear regression model, with the hypothesis h_{θ} , trained on dataset \mathbf{X} . This is expressed as:

$$\text{MSE}(\mathbf{X}, h_{\theta}) = J = \frac{1}{m} \sum_{i=1}^m (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})^2 \quad (2.3)$$

where the objective is to minimize the cost function by iteratively refining the values of $\boldsymbol{\theta}$.

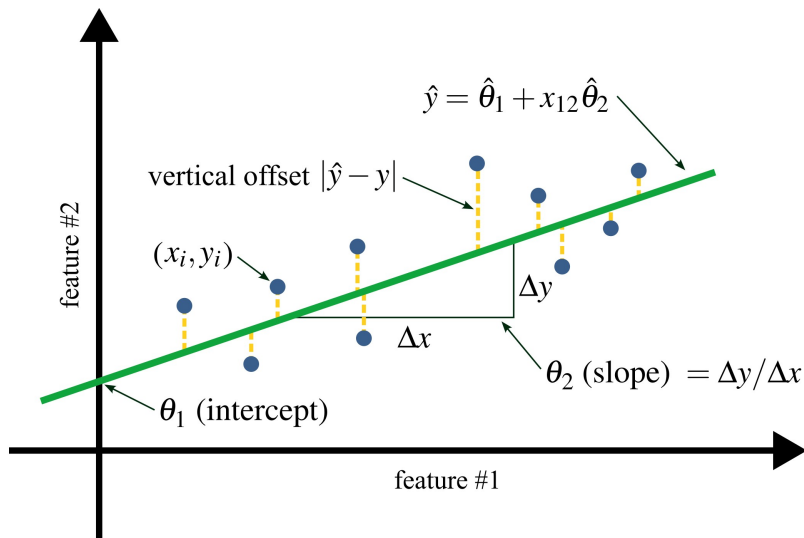


Figure 2.6: Linear regression of a two-feature dataset obtained with least squares.

Consider a dataset containing n observations, denoted as $(y_i, x_i)_{i=1}^n$. Each observation i consists of a scalar response y_i and a column vector x_i that holds the values for p predictors (regressors), represented as x_{ij} where $j = 1, \dots, p$. In the context of a linear regression model, the response variable y_i is modeled as a linear combination of the regressors, influenced by an error term ε_i :

$$y_i = \theta_1 x_{i1} + \theta_2 x_{i2} + \dots + \theta_m x_{im} + \varepsilon_i \quad (2.4)$$

This model can also be written in matrix notation as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\varepsilon} \quad (2.5)$$

Where:

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ x_{21} & x_{22} & \dots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{bmatrix}, \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad (2.6)$$

However, the best-fit model \hat{y} will not be able to account for the unmodeled error, therefore:

$$\hat{y}_i = \hat{\theta}_1 x_{i1} + \hat{\theta}_2 x_{i2} + \dots + \hat{\theta}_m x_{im} \quad (2.7)$$

or

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\theta}} \quad (2.8)$$

This goal is to find $\hat{\boldsymbol{\theta}}$ such that the error is minimized. Mathematically, this is simple enough as the $\hat{\boldsymbol{\theta}}$ is the minimization of the least-squares hyperplane, or:

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (2.9)$$

Equation 2.9 is referred to as the normal equation. Consider a simple linear model with $p = 2$, where $x_{11} = 1$. In this scenario, x_{11} acts as the bias term of the equation, whereas the remaining

elements of \mathbf{X} represent the data. Without this bias term, the solution would only address the slope of the line. Subsequently, solving for $\hat{\theta}_1$ and $\hat{\theta}_2$ yields:

$$[x_{11} x_{12}] [\hat{\theta}_1 \hat{\theta}_2] = \hat{y} \quad (2.10)$$

As $x_{11} = 1$ (again, called the bias term) this becomes:

$$\hat{y} = \hat{\theta}_1 + x_{12} \hat{\theta}_2 \quad (2.11)$$

This is simply another expression for the equation of a liner line:

$$y = mx + b \quad (2.12)$$

Solving the model output for a input or a series of inputs can than be done simply enough.

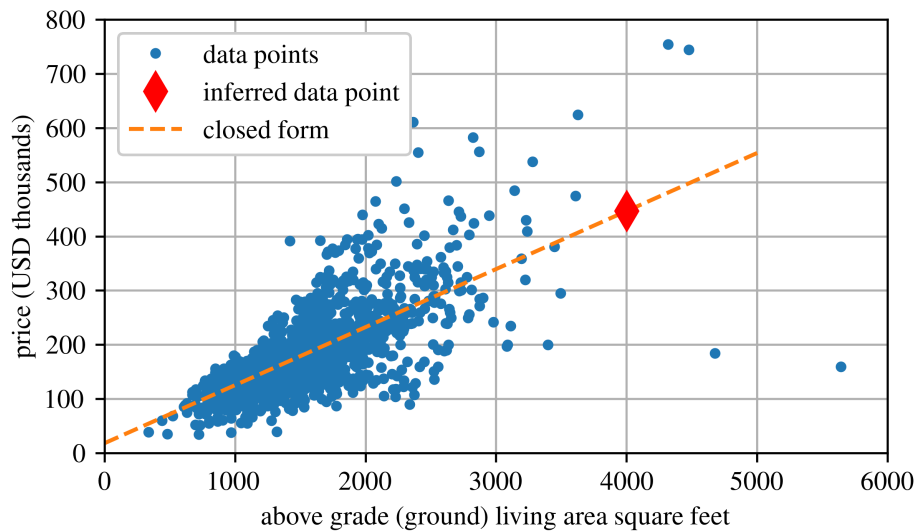


Figure 2.7: Ames housing data inferred.

2.1.2 Computational Complexity

The Normal Equation calculates the inverse of $\mathbf{X}^T \mathbf{X}$, an $n \times n$ matrix, where n represents the number of features. The computational complexity of matrix inversion generally ranges from $O(n^{2.4})$ to $O(n^3)$, depending on the specific algorithm used. Consequently, doubling the number of features would increase the computational effort by approximately $2^{2.4} \approx 5.3$ to $2^3 = 8$ times.

WARNING

When the number of features becomes very large, for example around 100,000, solving the Normal Equation becomes extremely slow.

On the positive side, the computational complexity of this equation is linear with respect to the number of instances in the training set, denoted as $O(m)$. This allows it to efficiently handle large training sets, as long as they fit within memory.

Furthermore, once your Linear Regression model is trained (whether through the Normal Equation or another method), making predictions is very quick. The computational effort is linear in relation to both the number of instances you wish to predict and the number of features. Essentially, predicting for twice as many instances or features will approximately double the computation time.

Next, we will explore different methods for training a Linear Regression model that are more appropriate for scenarios with a large number of features, or when the training data exceeds memory capacity.

Example 2.1 Linear Regression - Ames Housing Dataset

This example uses the Ames housing dataset to model the relationship between above-ground living area and house price using linear regression. It compares a manual fit, a closed-form solution, and gradient descent, alongside Scikit-Learn's `LinearRegression` and `SGDRegressor`, to illustrate how each method fits a line to the data.

Gradient Descent is a versatile optimization algorithm capable of finding optimal solutions across a broad spectrum of problems. The core concept of Gradient Descent involves iteratively adjusting parameters to minimize a cost function^a. Consider the analogy of being lost in a dense fog in the mountains and seeking to descend:

1. You sense the ground's slope under your feet (the local gradient).
2. You descend in the direction of the steepest descent.
3. Reaching a point where the local gradient is zero indicates that you have found a minimum.

The Gradient Descent process starts with random initialization of the parameter θ , followed by gradual improvements. Each incremental step aims to reduce the cost function until the algorithm converges to a minimum.

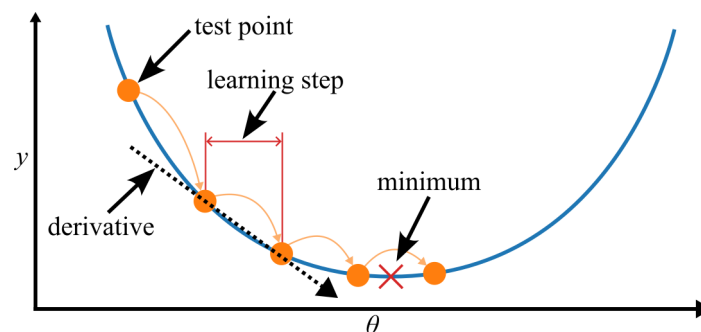


Figure 2.8: Illustration of gradient descent.

A critical aspect of Gradient Descent is the step size, controlled by the learning rate hyperparameter. A small learning rate leads to a slow convergence, requiring many iterations. Conversely, a high learning rate might cause the algorithm to overshoot the minimum, potentially causing divergence and failing to find an optimal solution.

^aRuder, Sebastian. "An overview of gradient descent optimization algorithms." arXiv preprint arXiv:1609.04747 (2016).

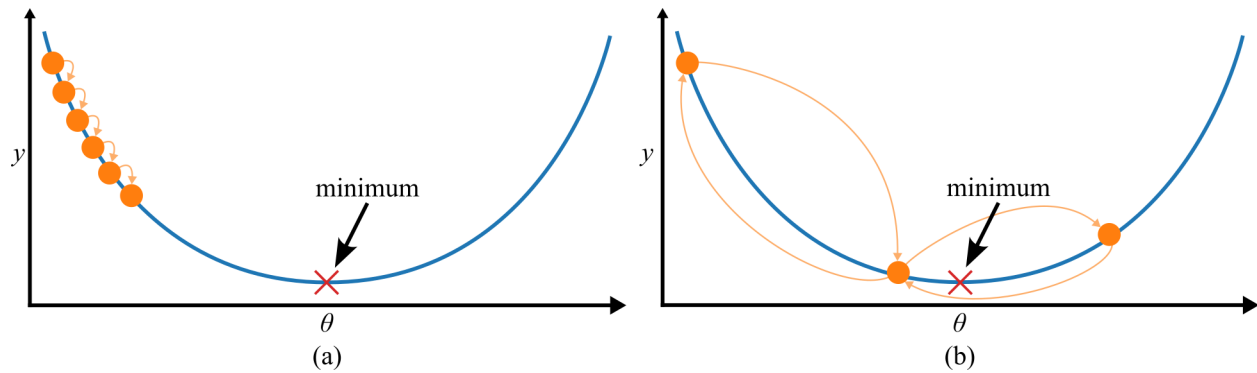


Figure 2.9: Effect of learning rate in gradient descent, showing: (a) a slow search with a low step size, and (b) a search with a large step size that oscillates around the minimum without ever finding the minimum.

Cost functions do not always present themselves as neat, concave shapes. They can feature obstacles such as holes, ridges, plateaus, and various complex topographies that complicate the convergence to the minimum. The challenges associated with Gradient Descent include:

- Starting the algorithm from a random point on the left may lead to convergence at a local minimum rather than the more optimal global minimum.
- Initiating on the right might result in a prolonged journey across a plateau, and terminating the process too soon could prevent reaching the global minimum.

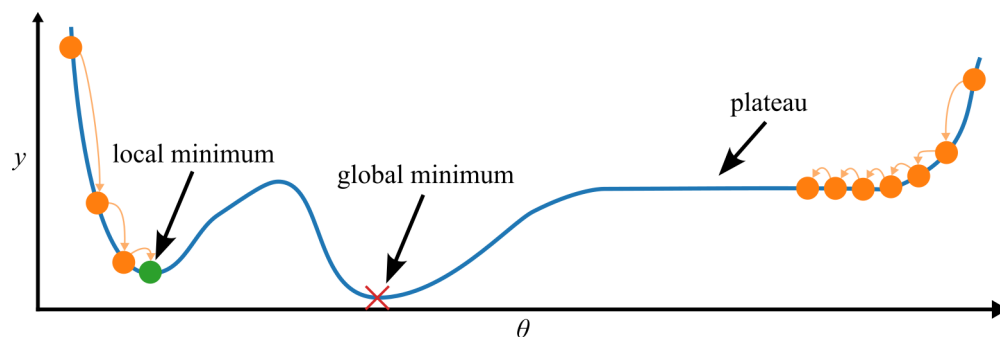


Figure 2.10: Gradient features and descent challenges.

Fortunately, the Mean Squared Error (MSE) cost function for Linear Regression models is convex. This characteristic ensures that for any two points on the curve, the line segment joining them does not intersect the curve itself. Consequently:

- There are no local minima, only one global minimum exists.
- It is a continuous function, and its slope changes smoothly.

These properties significantly benefit the optimization process: Gradient Descent can reliably approximate the global minimum given sufficient time and an appropriate learning rate.

2.1.3 Comparison of Gradient Descent Methods

Before looking into the mathematical details, let us examine how various gradient descent methods achieve the optimal solution, as illustrated in Figure 2.11. Additionally, Table 1 provides a comparative analysis of different gradient descent methods, highlighting their respective strengths and weaknesses.

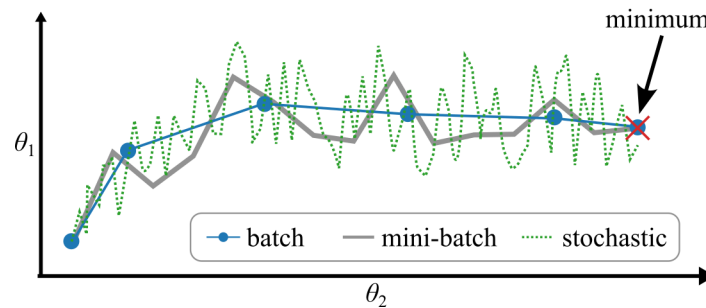


Figure 2.11: Comparing gradient descent methods.

Table 1: Comparison of various linear regression solutions.

algorithm	large number of instances (m)	data must fit in memory	large number of features (n)	hyperparams	scaling required	Scikit-Learn
normal equation	fast	yes	slow	0	no	LinearRegression
batch GD	slow	yes	fast	2	yes	n/a
stochastic GD	fast	no	fast	≥ 2	yes	SGDRegressor
mini-batch GD	fast	no	fast	≥ 2	yes	n/a

NOTE

After training, there is minimal difference between the models: these algorithms converge on similar solutions and make predictions in a nearly identical manner.

2.1.4 Batch Gradient Descent

Batch Gradient Descent (shown in figure 2.12) processes the entire batch of training data \mathbf{X} at every step, which is the origin of its name, “batch.” Consider again our Mean Squared Error (MSE) cost function:

$$J = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \quad (2.13)$$

For a hypothesis h_{θ} , the MSE for the dataset \mathbf{X} can be computed for each instance $\mathbf{x}^{(i)}$ as follows:

$$J(\mathbf{X}, h_{\theta}) = J = \frac{1}{m} \sum_{i=1}^m (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)})^2 \quad (2.14)$$

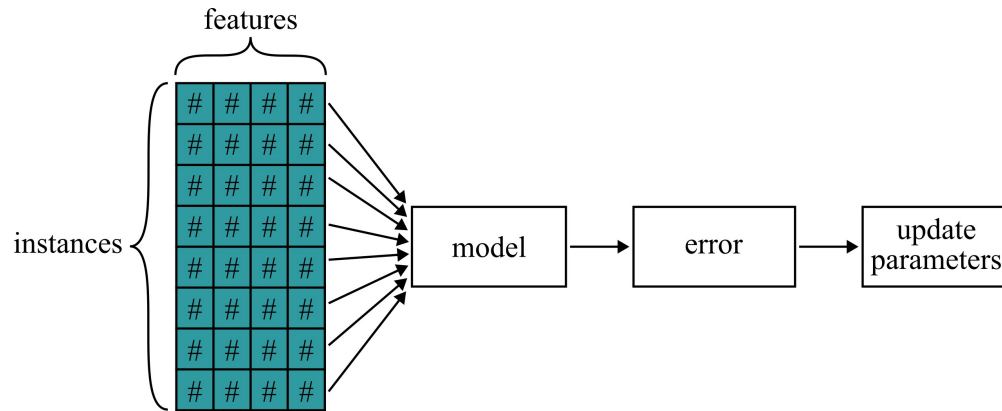


Figure 2.12: Flowchart of Batch Gradient Descent.

With this foundation, we can implement gradient descent by computing the gradient of the cost function with respect to each parameter θ_j . Specifically, this involves calculating how much the cost function changes when θ_j is altered slightly. This calculation is known as a partial derivative. Conceptually, it is akin to determining “the slope of the mountain under my feet if I face east,” then repeating the question facing north, and similarly for any other direction in a higher-dimensional space. The partial derivative of the cost function with respect to parameter θ_j is computed as:

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{2}{m} \sum_{i=1}^m (\boldsymbol{\theta}^T \mathbf{x}^{(i)} - y^{(i)}) x_j^{(i)} \quad (2.15)$$

Rather than computing each partial derivative individually, all partial derivatives can be calculated simultaneously using the gradient vector, $\nabla_{\theta} J(\boldsymbol{\theta})$, which encompasses all partial derivatives of the cost function for each model parameter:

$$\nabla_{\theta} J(\boldsymbol{\theta}) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} J(\boldsymbol{\theta}) \\ \frac{\partial}{\partial \theta_1} J(\boldsymbol{\theta}) \\ \vdots \\ \frac{\partial}{\partial \theta_n} J(\boldsymbol{\theta}) \end{pmatrix} = \frac{2}{m} \mathbf{X}^T \cdot (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) \quad (2.16)$$

WARNING

Be aware that this method computes using the entire training set \mathbf{X} at every step of Gradient Descent, which is why the approach is named Batch Gradient Descent. Although this method can be exceedingly slow with large training sets, it performs well with numerous features, outpacing the Normal Equation in training Linear Regression models with extensive features.

Once you obtain the gradient vector, which indicates the direction of ascent, the next step involves moving in the opposite direction, effectively going downhill. This is achieved by subtracting $\nabla_{\theta} J(\theta)$ from θ , incorporating the learning rate η to scale the step size:

$$\theta^{(\text{next step})} = \theta - \eta \nabla_{\theta} J(\theta) \quad (2.17)$$

Interestingly, this process aligns perfectly with the results from the Normal Equation. However, variations in the learning rate η can significantly influence the outcomes. Figure 2.13 illustrates the first 10 steps of Gradient Descent with three different learning rates, with the dashed line marking the starting point.

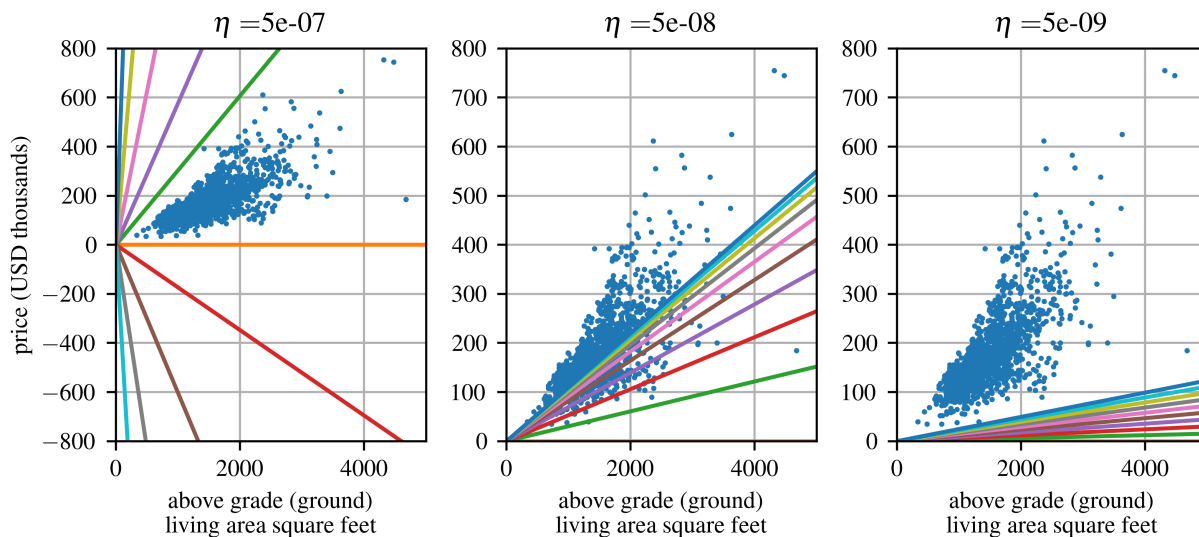


Figure 2.13: The effects of different learning rates on the gradient descent algorithms, showing: (a) too high ($\eta = 5e - 07$); (b) about right ($\eta = 5e - 08$), and; (c) too low ($\eta = 5e - 09$).

Figure 2.13 highlights the impact of the learning rate on gradient descent. Figure 2.13(a) With a very large step size, $\eta = 5 \times 10^{-7}$, each update overshoots the minimum so the algorithm diverges rather than converging. Figure 2.13(b) Using $\eta = 5 \times 10^{-8}$ provides a balanced step size; the cost decreases just right^a and the algorithm reaches the optimum in only a few iterations. Figure 2.13(c) A small step size, $\eta = 5 \times 10^{-9}$, keeps the algorithm moving toward the minimum yet progress is slow, requiring many more iterations to achieve the same result.

To identify an appropriate learning rate, a grid search can be employed. It is advisable to limit the number of iterations during grid search to avoid models that are too slow to converge.

Determining the correct number of iterations is also crucial. If set too low, the algorithm may stop far from the optimal solution. Conversely, overly high settings lead to unnecessary computations after convergence. A practical approach is to allow a large number of iterations but to halt the algorithm when the gradient vector's norm shrinks to below a small threshold ϵ (known as the tolerance), indicating proximity to the minimum.

^aPyle, K. "Goldilocks and the three bears." *Mother's Nursery Tales* (1918): 207-213.

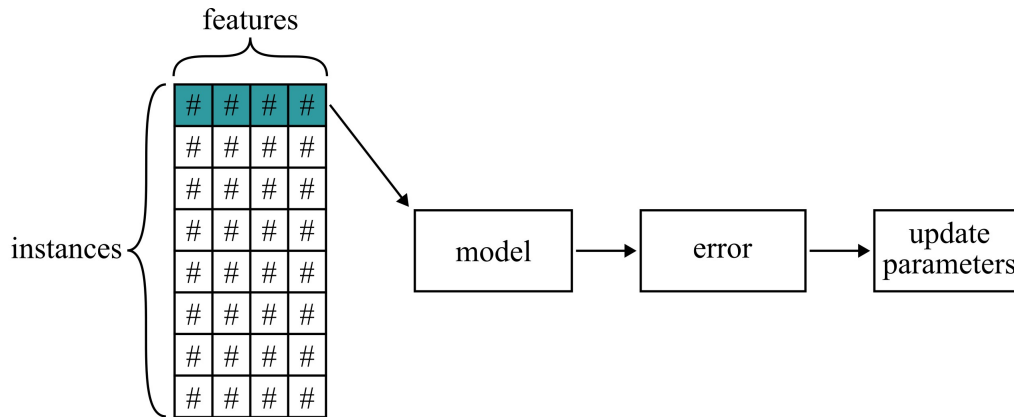


Figure 2.14: Flowchart of Stochastic Gradient Descent.

2.1.5 Stochastic Gradient Descent

Stochastic Gradient Descent (shown in figure 2.14) updates the model parameters after evaluating each individual training sample $x^{(i)}$ and its corresponding label $y^{(i)}$, which stands in contrast to batch gradient descent that uses the entire dataset for each update. The update rule for SGD can be expressed as

$$\theta^{(\text{next step})} = \theta - \eta \nabla_{\theta} J(\theta; x^{(i)}; y^{(i)}). \quad (2.18)$$

While batch gradient descent tends to perform unnecessary recalculations for large datasets by reevaluating gradients for similar examples with each update, SGD eliminates this inefficiency by updating parameters incrementally. Consequently, SGD is generally faster and adaptable for online learning. Due to its frequent and individual updates, SGD exhibits high variance in the objective function, causing significant fluctuations as illustrated in Fig. 2.15(b).

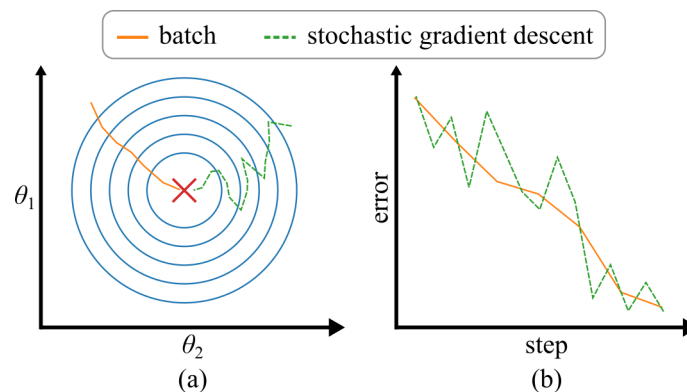


Figure 2.15: Comparison of the iterative performance of Batch and Stochastic Gradient Descent, showing: (a) how they move towards their target, and (b) the error at comparable steps.

2.1.6 Mini-batch Gradient Descent

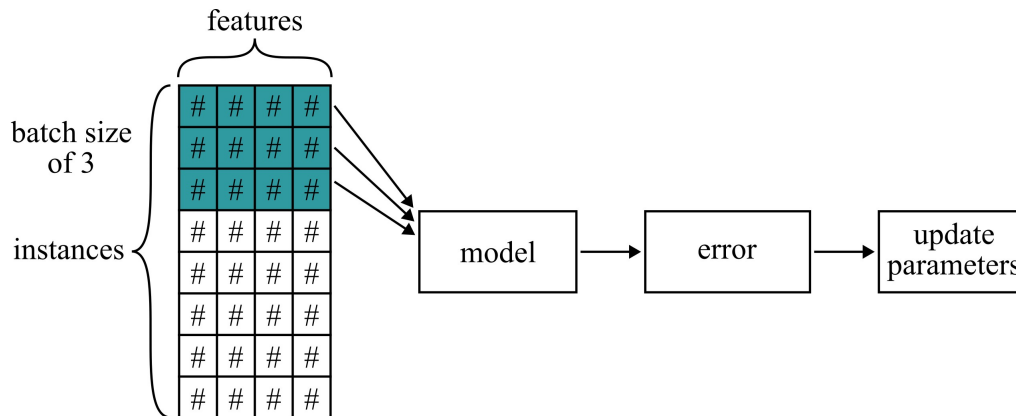


Figure 2.16: Flowchart of Mini-batch Gradient Descent.

$$\theta^{(\text{next step})} = \theta - \eta \nabla_{\theta} J(\theta; x^{(i:i+n)}; y^{(i:i+n)}) \quad (2.19)$$

The final Gradient Descent algorithm we will discuss is Mini-batch Gradient Descent. This method combines elements of both Batch and Stochastic Gradient Descent: rather than calculating gradients using the entire training set (as in Batch GD) or a single instance (as in Stochastic GD), Mini-batch GD computes gradients using small, randomly selected subsets of instances known as minibatches. One significant benefit of Mini-batch GD over Stochastic GD is the improved performance from hardware-optimized matrix operations, particularly on GPUs.

Mini-batch GD's trajectory through parameter space tends to be more stable compared to the often erratic path of SGD, especially with larger minibatches. This stability typically brings Mini-batch GD closer to the minimum than SGD. However, Mini-batch GD might struggle more with escaping local minima in scenarios prone to such issues, which is less of a concern in Linear Regression. Figure 2.11 illustrates the trajectories of these three Gradient Descent techniques during training. While Batch GD precisely reaches the minimum, Stochastic GD and Mini-batch GD tend to oscillate nearby. It is important to note, however, that while Batch GD is slow in taking steps, both Stochastic GD and Mini-batch GD could also effectively reach the minimum with an appropriate learning rate schedule.

2.1.7 Feature Scaling

The cost function can resemble an elongated bowl if the feature scales vary significantly. Figure 2.17 illustrates the effect of feature scaling on Gradient Descent. In figure 2.17(a), Gradient Descent progresses directly towards the minimum, reaching it swiftly. However, in figure 2.17(b) it initially moves almost orthogonally to the direction of the global minimum, followed by a lengthy traverse down a nearly flat valley, significantly delaying the convergence.

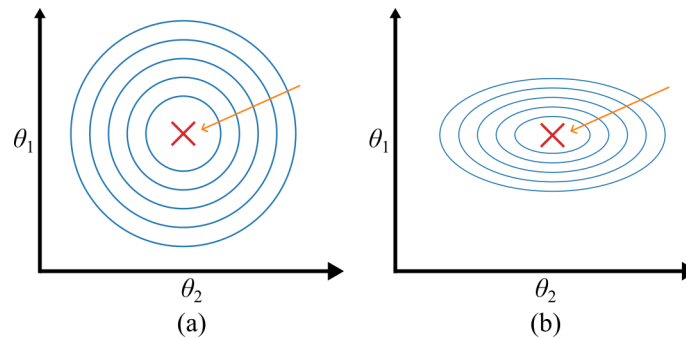


Figure 2.17: Gradient descent on a 2D surface for parameters that are: (a) equally scaled, and (b) unequally scaled.

Training a model can be viewed as a search through the parameter space for the combination that minimizes the cost function. When additional parameters are introduced, this space gains extra dimensions, making the optimization task more challenging. In Ordinary Least Squares Linear Regression, however, the cost surface is convex and bowl-shaped, so any descent method is guaranteed to reach the global minimum. To normalize feature scales, two common methods are employed. They are Min-max scaling and Standardization.

- **Min-max scaling**, often referred to as normalization, is straightforward: it rescales the data to a range of 0 to 1 by subtracting the minimum value and dividing by the range (max minus min).

$$\mathbf{X}' = \frac{\mathbf{X} - \mathbf{X}_{\min}}{\mathbf{X}_{\max} - \mathbf{X}_{\min}} \quad (2.20)$$

Scikit-Learn offers `sklearn.preprocessing.MinMaxScaler` for this purpose.

- **Standardization** differs significantly as it first subtracts the mean (resulting in a zero mean) and then divides by the standard deviation to achieve unit variance. This method does not limit values to a specific range, which can be problematic for certain algorithms, such as neural networks which often expect inputs between 0 and 1. However, standardization is less sensitive to outliers. For instance, if a median income is mistakenly recorded as 100, min-max scaling would compress all other values between 0 and 15 to between 0 and 0.15, whereas standardization would be minimally affected.

$$\mathbf{X}' = \frac{\mathbf{X} - \bar{\mathbf{X}}}{\sigma} \quad (2.21)$$

`sklearn.preprocessing.StandardScaler` is provided by Scikit-Learn for standardization.

WARNING

It is crucial to apply scalers exclusively to the training data and not to the entire dataset, which includes the test set. This practice ensures that the model is not inadvertently exposed to test data during training. Once the scalers are fitted to the training data, they can then be used to transform the training set, the test set, and any new data subsequently encountered.

2.2 Polynomial Regression

What if the underlying pattern of your data is more complex than a simple linear relationship? Consider the dataset shown in figure 2.18 with non-linear characteristics (one feature across multiple instances), which typically can be modeled using a second-order polynomial

$$\hat{y} = ax^2 + bx + c. \quad (2.22)$$

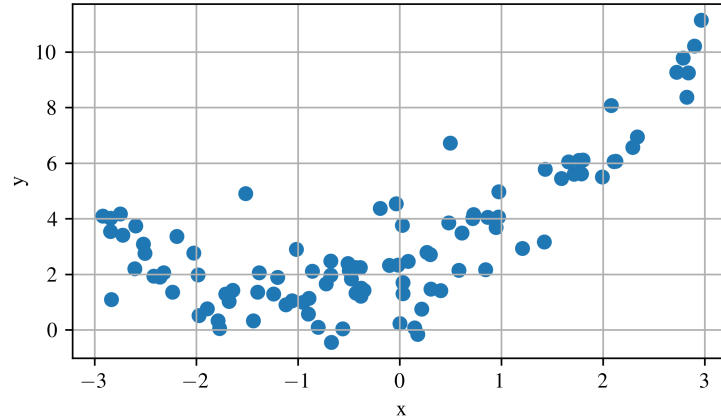


Figure 2.18: Dataset derived from a 2nd-order polynomial.

A simple linear fit will under-perform, so you first enrich the training set by adding the squared term of each feature as an extra column; shown in figure 2.19. Although the model you train is still linear with respect to its parameters, it now operates on an augmented feature space and can represent the desired quadratic curve. This approach is called Polynomial Regression, and it extends naturally by including higher powers whenever more complex nonlinearities are present.

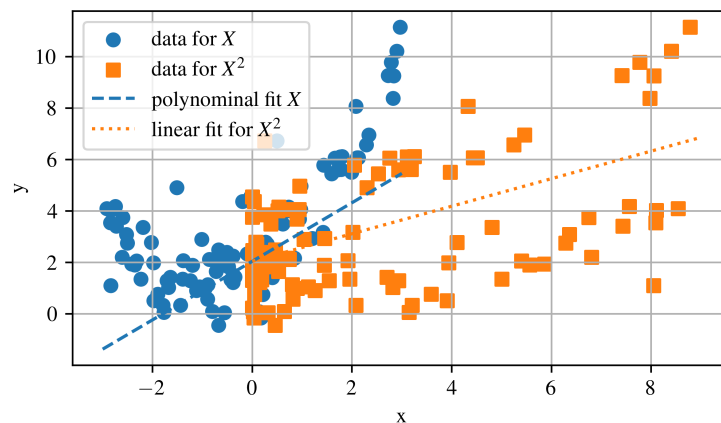


Figure 2.19: Linear models fit to the individual features of a polynomial dataset.

This data can then be incorporated into two linear models, where the slopes of the feature sets serve as the parameters for the base-line polynomial expression (a and b in Equation 2.22) and the bias term is the offset (c in Equation 2.22). The results of such a polynomial fit is shown in figure 2.20.

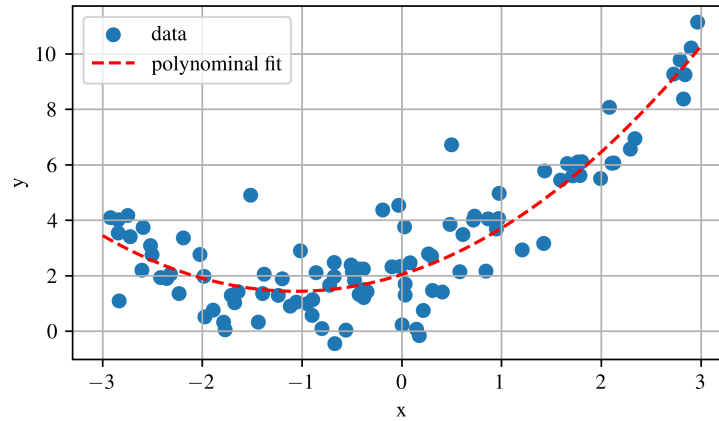


Figure 2.20: Polynomial model fit to a dataset.

It is important to note that Polynomial Regression can identify interrelationships between features in cases where multiple features exist, a capability beyond the scope of simple Linear Regression. This enhancement is enabled by `PolynomialFeatures`, which includes all possible combinations of features up to the specified degree. For instance, with two features a and b , and a degree of 3, `PolynomialFeatures` would add not only a^2 , a^3 , b^2 , and b^3 but also the combined terms ab , a^2b , and ab^2 .

WARNING

`PolynomialFeatures(degree=d)` expands an array with n original features into one that includes $\frac{(n+d)!}{d!n!}$ features, accounting for all combinations of features up to the d -th degree. Here, $n!$ represents the factorial of n , calculated as $1 \times 2 \times 3 \times \dots \times n$. Be cautious of the rapid increase in the number of features, known as combinatorial explosion!

Example 2.2 Polynomial Regression

This example fits a non-linear dataset using polynomial regression by adding x^2 as an extra feature. It demonstrates how linear models can approximate curved relationships when polynomial terms are included, and shows the resulting fit compared to the original data.

2.3 Training and Testing Data

Training a predictive model on a dataset and then testing it with the same data is a fundamental error in methodology. Such a model could simply memorize the labels of the training samples, achieving perfect performance during training but failing to make any useful predictions on new, unseen data. This phenomenon is known as overfitting. To prevent this, it is standard practice in supervised machine learning to reserve a portion of the available data as a test set, denoted as \mathbf{X}_{test} and \mathbf{y}_{test} .

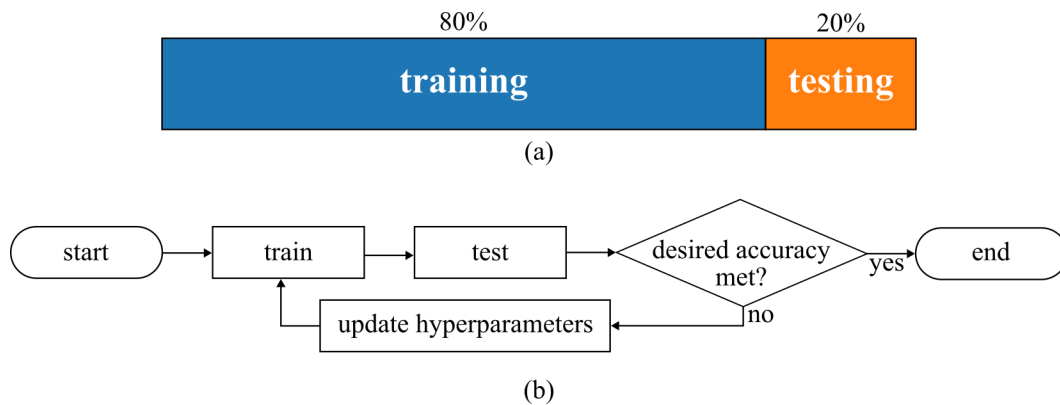


Figure 2.21: Splitting data up into training and testing subsets.

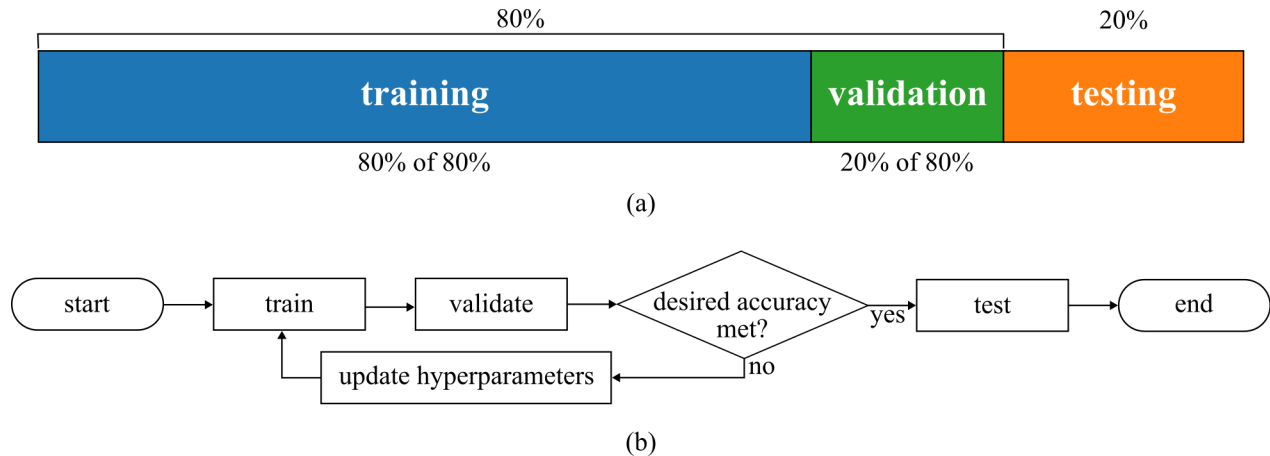


Figure 2.22: Splitting data up into training, validation, and testing subsets.

The function `sklearn.model_selection.train_test_split` in scikit-learn randomly divides arrays or matrices into training and testing subsets.

2.4 Pipelines

Pipelines in scikit-learn streamline the process by sequentially applying a list of transformations followed by a final estimator to a dataset. Employing pipelines allows for the integration of multiple processing steps, which can then be cross-validated together while experimenting with various parameters. Fig. 2.23 illustrates a typical pipeline configuration in scikit-learn.

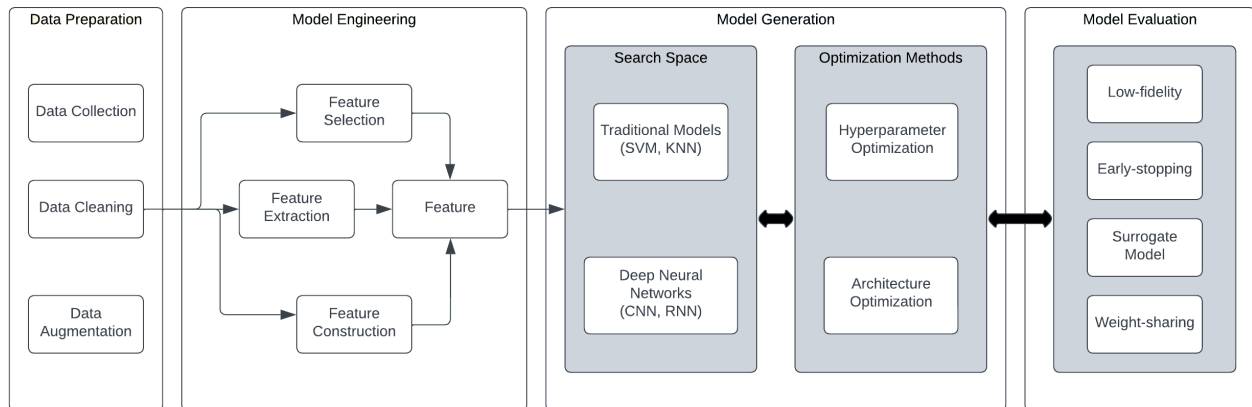


Figure 2.23: Pipeline setup for the automated deployment of pre-processing and modeling steps.^a

2.5 Learning Curves

Performing Polynomial Regression with a high degree typically allows for a closer fit to the training data compared to simple Linear Regression. For instance, Figure 2.24 demonstrates the application of a 30-degree polynomial model to a set of training data, contrasting it with both a linear model and a quadratic model (2nd-degree polynomial). Observe how the 30-degree polynomial model contorts to closely match the training instances.

- The linear model exhibits underfitting.
- The high-degree Polynomial Regression model drastically overfits the training data.
- Among these, the quadratic model is likely to generalize best.

In this instance, the appropriateness of the quadratic model is clear since the data was initially generated using such a model. However, in real-world scenarios where the underlying function of the data is unknown, determining the optimal complexity for your model can be challenging. How can you ascertain whether your model is overfitting or underfitting the data?

^aModified from PopovaZhuhadar, CC BY-SA 4.0 <<https://creativecommons.org/licenses/by-sa/4.0/>>, via Wikimedia Commons

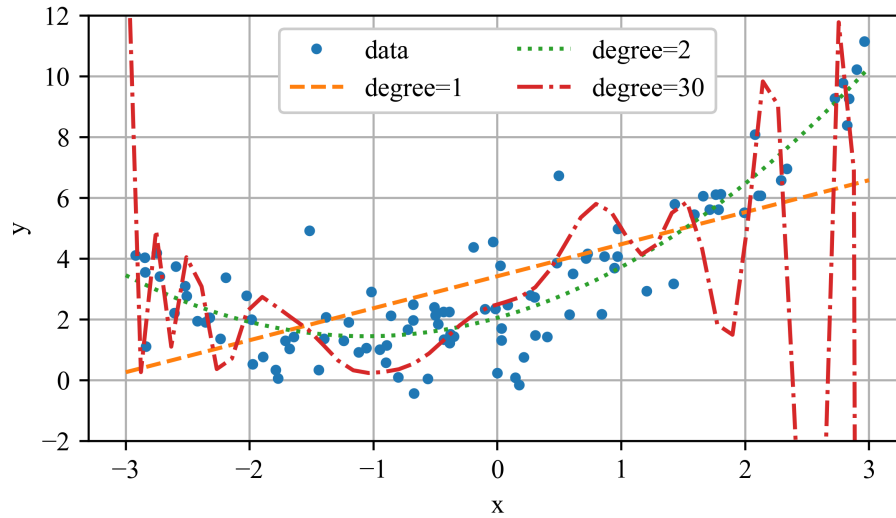


Figure 2.24: Polynomial regression showing underfitting (degree=1), a respectable model fit (degree=2), and overfitting (degree=30).

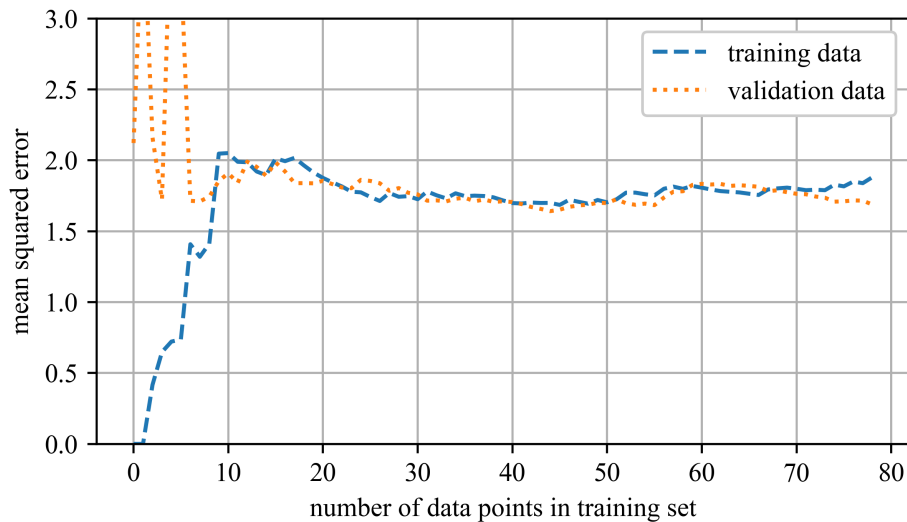


Figure 2.25: Learning curves for the underfitting linear model (degree=1) where underfitting is evident because both curves have plateaued; they are close together and relatively high.

Examining the model's behavior in figures 2.25 and 2.26 on the training set shows a clear trend. With only one or two samples the fit is perfect, so the error curve begins at zero. As additional observations are introduced the model can no longer capture every point exactly, partly because of measurement noise and partly because the underlying relationship is nonlinear, and consequently the training error rises. After a sufficient number of samples the curve flattens out; beyond this plateau adding further data neither markedly lowers nor raises the average error.

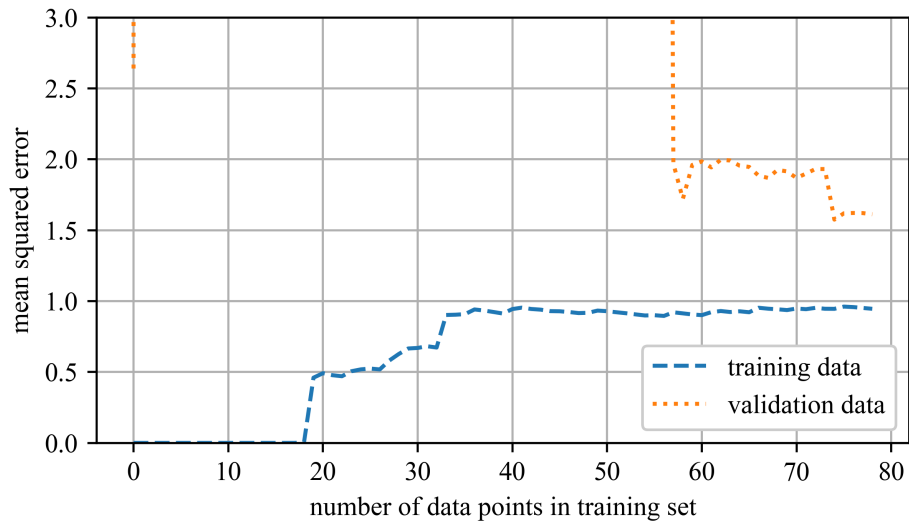


Figure 2.26: Learning curves for the overfitting polynomial model (degree=20) showing a significant gap between the curves, which indicates better performance on the training data than on the validation data.

Inspection of the validation curve in figure 2.26 tells a complementary story. With only a handful of training samples the model cannot generalize, so the validation error starts high. As more examples become available it learns progressively better patterns and the error falls. Eventually, though, the simple linear hypothesis is unable to capture the data's full complexity, causing the decline to flatten out; the validation curve plateaus at nearly the same level as the training curve. Their close proximity and uniformly large values are characteristic of an underfitting model.

Figure 2.26 displays the learning curves for a model fitted with a 20th-degree polynomial. The overall shape resembles the curves seen earlier, yet two key differences stand out. First, the error on the training set is far lower than that achieved by the Linear Regression model, demonstrating the polynomial's extra flexibility. Second, a pronounced gap separates the training and validation curves, which means the model performs much better on the data it has already seen; this divergence is the hallmark of overfitting. Expanding the training set could reduce the gap, but without additional regularization the risk of overfitting would remain.

Next, let's apply the code using a 2nd order polynomial, which accurately captures the essence of the data without underfitting or overfitting. These results are shown in Figure 2.27. Here it can be seen that the training and validation curves again meet after converging, showing a well-fit model.

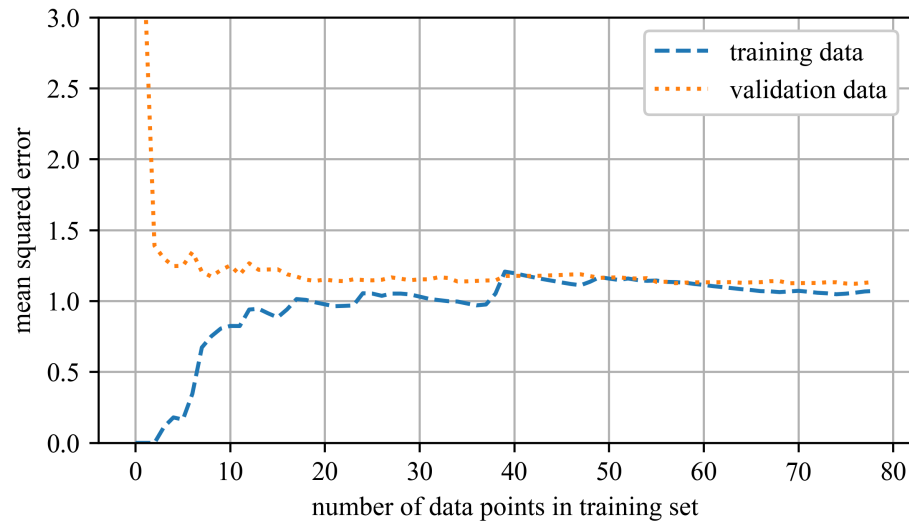


Figure 2.27: Learning curves for the optimal polynomial model (degree=2) showing well-aligned curves that plateau at a relatively low error level.

Example 2.3 Learning Curves

This example compares learning curves for linear and polynomial regression models. By plotting training and validation errors as the training set grows, it illustrates how model complexity impacts bias and variance, and helps identify underfitting and overfitting behavior.

2.6 Regularized Linear Models

Overfitting is frequently curbed by regularizing the model, meaning you introduce constraints that reduce its effective degrees of freedom and make it harder to memorize noise. In a polynomial setting you can achieve this simply by choosing a lower-degree polynomial, while in a linear model you typically add a penalty that discourages large coefficient values. This naturally raises the question: how is regularization expressed mathematically?

2.6.1 Ridge Regression

Ridge Regression modifies Linear Regression by adding a regularization term to the cost function, which compels the learning algorithm to fit the data while keeping the model weights as small as possible. The updated cost-function is

$$J(\theta) = \text{MSE}(\theta) + \alpha \frac{1}{2} \sum_{i=1}^n \theta_i^2. \quad (2.23)$$

The hyperparameter α dictates the degree of regularization. If $\alpha = 0$, Ridge Regression reduces to plain Linear Regression. If α is very high, the model weights are significantly shrunk, resulting in a nearly flat line through the mean of the data. Equation 2.23 outlines the cost function for Ridge Regression. The bias term θ_0 is not regularized, as indicated by the summation starting from $i = 1$.

NOTE

The regularization term is only included during the training phase. After training, the model's performance should be evaluated based on the unregularized performance measure.

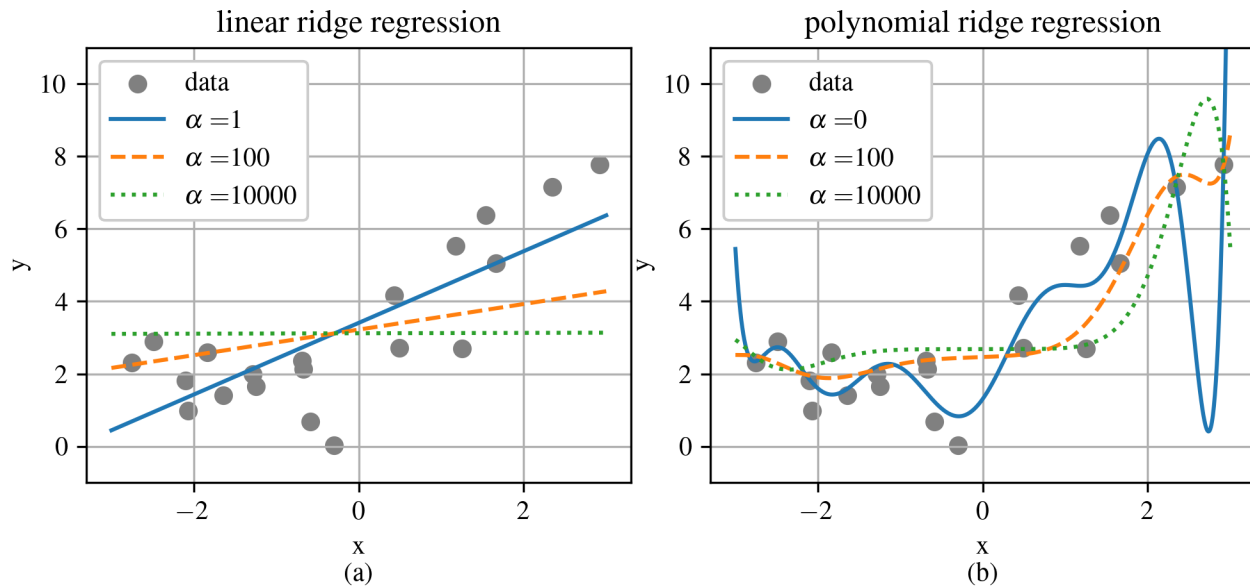


Figure 2.28: Visualization of Ridge Regression models with varying α values.

Figure 2.28 illustrates various Ridge models trained on linear data, showcasing different α value. In Figure 2.28(a), models are purely Ridge Regression, leading to linear predictions. In Figure 2.28(b), the data undergoes Polynomial Regression with Ridge regularization. To do this, the data is first expanded using `PolynomialFeatures` (`degree=10`), then scaled with `StandardScaler`, and finally, Ridge models are applied to these transformed features.

Increasing α results in flatter (i.e., more reasonable and less extreme) predictions, thereby reducing the model's variance but increasing its bias. Ridge Regression can be performed using a closed-form solution or through Gradient Descent, similar to Linear Regression. The pros and cons for each method mirror those discussed previously. The closed-form solution, an extension of the normal equation, is given by:

$$\hat{\theta} = (X^T \cdot X + \alpha A)^{-1} \cdot X^T \cdot y \quad (2.24)$$

where α influences the regularization and A is an $n \times n$ identity matrix with the top-left value replaced by 0 to exclude the bias term. When using Gradient Descent, the derivative of the cost function guides the adjustment toward optimal model parameters.

Example 2.4 Ridge Regression

This example demonstrates Ridge Regression as a regularized alternative to linear regression. A high-degree polynomial model is fit to noisy data using a pipeline with Ridge regularization, showing how the regularization term reduces overfitting and stabilizes the model.

2.7 Early Stopping

An alternative approach to regularizing iterative learning methods (such as Gradient Descent) is to halt training as soon as the validation error hits its lowest point, a strategy known as early stopping. Figure 2.29 illustrates a complex model, specifically a high-degree Polynomial Regression model, being refined through Batch Gradient Descent. Over the course of training, as epochs progress, the model's prediction error (RMSE) on both the training and validation sets decreases. However, after some time, the validation error ceases to decline and begins to increase, signaling that the model is beginning to overfit the training data. By implementing early stopping, training is terminated the moment the validation error reaches its minimum. Geoffrey Hinton praised this method stating, "Early stopping (is) beautiful free lunch" due to its simplicity ^a.

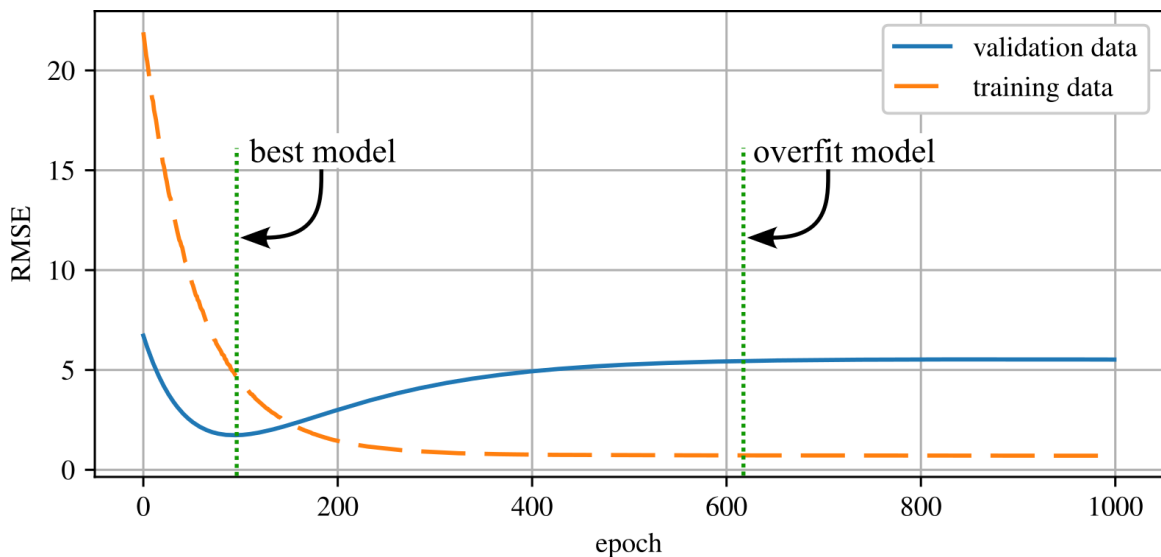


Figure 2.29: An example showing how early stopping can result in a better model as it does not allow for the over-fitting of the training set.

Example 2.5 Early Stopping

This example uses a high-degree polynomial model trained with stochastic gradient descent to demonstrate early stopping. By tracking training and validation errors across epochs, it highlights how halting training early can prevent overfitting and improve generalization.

^aPennington, Jeffrey, Richard Socher, and Christopher D. Manning. "Glove: Global vectors for word representation." Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). 2014.

2.8 Examples

Example 2.1

```

1  """
2  Example 2.1 Linear Regression
3  @author: Austin Downey
4  """
5
6  import IPython as IP
7  IP.get_ipython().run_line_magic('reset', '-sf')
8
9  import numpy as np
10 import matplotlib.pyplot as plt
11 import sklearn as sk
12 from sklearn import datasets
13 from sklearn.linear_model import LinearRegression
14
15 plt.close('all')
16
17 """ load data
18
19 ames = sk.datasets.fetch_openml(name="house_prices", as_frame=True, parser='auto')
20 target = ames['target'].values
21 data = ames['data']
22 YrSold = data['YrSold'].values # Year Sold (YYYY)
23 MoSold = data['MoSold'].values # Month Sold (MM)
24 OverallCond = data['OverallCond'].values # OverallCond: Rates the overall condition of
the house
25 GrLivArea = data['GrLivArea'].values # Above grade (ground) living area square feet
26 BedroomAbvGr = data['BedroomAbvGr'].values # Bedrooms above grade (does NOT include
basement bedrooms)
27
28 # Ask a home buyer to describe their dream house, and they probably won't begin
29 # with the height of the basement ceiling or the proximity to an east-west railroad.
30 # But this playground competition's dataset proves that much more influences price
31 # negotiations than the number of bedrooms or a white-picket fence.
32
33 # With 79 explanatory variables describing (almost) every aspect of residential
34 # homes in Ames, Iowa, this competition challenges you to predict the final price
35 # of each home.
36
37 # Plot a few of the interesting features vs the target (price). In particular,
38 # let's plot the number of rooms vs. the price.
39 plt.figure()
40 plt.plot(GrLivArea, target, 'o', markersize=2)
41 plt.xlabel('Above grade (ground) living area square feet')
42 plt.ylabel('price (USD)')
43 plt.grid(True)
44 plt.tight_layout()
45
46 """ Build a model for the data
47 X = GrLivArea
48 Y = target
49 model_X = np.linspace(0, 5000)
50
51 theta_1 = 0
52 theta_2 = 100
53 model_Y_manual = theta_1 + theta_2*model_X
54
55 plt.figure()
56 plt.plot(X, Y, 'o', markersize=2, label='data')
57 plt.plot(model_X, model_Y_manual, '--', label='manual fit')
58 plt.xlabel('Above grade (ground) living area square feet')
59 plt.ylabel('price (USD)')
60 plt.grid(True)
61 #plt.xlim([3.5, 9])
62 #plt.ylim([0, 50000])
63 plt.legend(framealpha=1)
64 plt.tight_layout()
65

```

```

66 # add a dimension to the data as math is easier in 2d arrays and sk learn only
67 # takes 2d arrays
68 X = np.expand_dims(X,axis=1)
69 Y = np.expand_dims(Y,axis=1)
70 model_X = np.expand_dims(model_X,axis=1)
71
72 ### compute the linear regression solution using the closed form solution
73
74 # compute
75 X_b = np.ones((X.shape[0],2))
76 X_b[:,1] = X.T # add x0 = 1 to each instance
77
78 theta_closed_form = np.linalg.inv(X_b.T@X_b)@X_b.T@Y
79
80 model_y_closed_form = theta_closed_form[0] + theta_closed_form[1]*3000
81 model_Y_closed_form = theta_closed_form[0] + theta_closed_form[1]*model_X
82
83 plt.figure()
84 plt.plot(X,Y,'o',markersize=3,label='data points')
85 plt.xlabel('Above grade (ground) living area square feet')
86 plt.ylabel('price (USD)')
87 plt.plot(3000,model_y_closed_form,'dr',markersize=10,zorder=10,
88         label='inferred data point')
89 plt.plot(model_X,model_Y_closed_form,'-',label='normal equation')
90 plt.grid(True)
91 plt.legend()
92 plt.tight_layout()
93
94 ### compute the linear regression solution using gradient descent
95
96 eta = 0.00000001 # learning rate
97 n_iterations = 100
98 m = X.shape[0]
99 theta_gradient_descent = np.random.randn(2,1) # random initialization
100 for iteration in range(n_iterations):
101     gradients = 2/m * X_b.T.dot(X_b.dot(theta_gradient_descent) - Y)
102     theta_gradient_descent = theta_gradient_descent - eta * gradients
103
104 print(theta_gradient_descent)
105
106 model_Y_gradient_descent = theta_gradient_descent[0] \
107     + theta_gradient_descent[1]*model_X
108
109 plt.figure()
110 plt.plot(X,Y,'o',markersize=3,label='data points')
111 plt.xlabel('Above grade (ground) living area square feet')
112 plt.ylabel('price (USD)')
113 plt.plot(model_X,model_Y_closed_form,'-',label='normal equation')
114 plt.plot(model_X,model_Y_gradient_descent,':',label='gradient descent')
115 plt.grid(True)
116 plt.legend()
117 plt.tight_layout()
118
119 ### compute the linear regression solution using sk-learn
120
121 # build and train a closed from linear regression model in sk-learn
122 model_LR = sk.linear_model.LinearRegression()
123 model_LR.fit(X,Y[:,0])
124 model_Y_sk_LR = model_LR.predict(model_X)
125
126 # build and train a Stochastic Gradient Descent linear regression model in sk-learn.
127 # Note that in running the model, the best way to do this would be to use a pipeline
128 # =with feature scaling. However, here we just set 'eta0' to a low value, this
129 # is done only for educational # purposes and is not the ideal methodology in
130 # terms of system robustness.
131 model_SGD = sk.linear_model.SGDRegressor(learning_rate='constant',eta0=0.00000001)
132 model_SGD.fit(X,Y[:,0])

```

```
133 model_Y_sk_SGD = model_SGD.predict(model_X)
134
135 # plot the modeled results
136 plt.figure()
137 plt.plot(X,Y,'o',markersize=2,label='data')
138 plt.plot(model_X,model_Y_closed_form,'-',label='normal equation')
139 plt.plot(model_X,model_Y_gradient_descent,'--',label='gradient descent')
140 plt.plot(model_X,model_Y_sk_LR,':',label='sklearn normal equation')
141 plt.plot(model_X,model_Y_sk_SGD,'-.',label='sklearn stochastic gradient descent')
142
143 plt.xlabel('Above grade (ground) living area square feet')
144 plt.ylabel('price (USD)')
145 plt.grid(True)
146 plt.legend(framealpha=1)
147 plt.tight_layout()
```

Example 2.2

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Example 2.2
5  Polynomial regression
6  Machine Learning for Engineering Problem Solving
7  @author: Austin Downey
8  """
9
10 import IPython as IP
11 IP.get_ipython().run_line_magic('reset', '-sf')
12
13 import numpy as np
14 import scipy as sp
15 import matplotlib as mpl
16 import matplotlib.pyplot as plt
17 import sklearn as sk
18 from sklearn import linear_model
19
20 plt.close('all')
21
22 """ build the data sets
23 np.random.seed(2) # 2 and 6 are pretty good
24 m = 100
25 X = 6 * np.random.rand(m,1) - 3
26 Y = 0.5 * X**2 + X + 2 + np.random.randn(m,1)
27
28 # plot the data
29 plt.figure()
30 plt.grid(True)
31 plt.scatter(X,Y)
32 plt.xlabel('x')
33 plt.ylabel('y')
34
35 """ perform polynomial regression
36
37 # generate x^2 as we use the model y = a*x^2 + b*x + c
38 X_poly_manual = np.hstack((X,X**2))
39
40 # or use the code as this does lots of features for multi-feature data sets.
41 poly_features = sk.preprocessing.PolynomialFeatures(degree=2, include_bias=False)
42 X_poly_sk = poly_features.fit_transform(X)
43
44 # they do do same thing as shown below, so select one to carry forward.
45 print(X_poly_manual == X_poly_sk)
46 X_poly = X_poly_manual
47
48 # In essence, we now have two data sets. We can plot that here
49 plt.figure()
50 plt.grid(True)
51 plt.scatter(X_poly[:,0],Y,label = 'data for x')
52 plt.scatter(X_poly[:,1],Y,marker='s',label = 'data for $x^2$')
53 plt.legend()
54 plt.xlabel('x')
55 plt.ylabel('y')
56
57 # and fit linear models to these data sets
58 model = sk.linear_model.LinearRegression() # Select a linear model
59 model.fit(X_poly,Y) # Train the model
60 X_model_1 = np.linspace(-3,3)
61 X_model_2 = np.linspace(0,9)
62
63 # the model parameters are:
64 model_coefficients = model.coef_
65 model_intercept = model.intercept_
66 print(model_coefficients)
67 print(model_intercept)

```

```
68
69 Y_X1 = model_coefficients[0][0]*X_model_1 + model_intercept
70 Y_X2 = model_coefficients[0][1]*X_model_2 + model_intercept
71
72 # now if we plot the linear models on the extended set of features.
73 plt.figure()
74 plt.grid(True)
75 plt.scatter(X_poly[:,0],Y,label = 'data for x')
76 plt.scatter(X_poly[:,1],Y,marker='s',label = 'data for $x^2$')
77 plt.plot(X_model_1,Y_X1,'--',label='linear fit $x$')
78 plt.plot(X_model_2,Y_X2,':',label='linear fit for $x^2$',)
79 plt.legend()
80 plt.xlabel('x')
81 plt.ylabel('y')
82 plt.savefig('example_6_fig_1',dpi=300)
83
84 # now that we have a parameter for x and x^2, these can be recombined into a single
85 # model, y = x^2*a + x*b + c.
86 Y_polynomial = model_coefficients[0][1]*X_model_1**2 + model_coefficients[0][0]*\
87     X_model_1 + model_intercept
88
89 plt.figure()
90 plt.grid(True)
91 plt.scatter(X,Y,label='data')
92 plt.plot(X_model_1,Y_polynomial,'r--',label='polynomial fit')
93 plt.xlabel('x')
94 plt.ylabel('y')
95 plt.legend()
96 plt.savefig('example_6_fig_2',dpi=300)
97
98
99
100
101
```

Example 2.3

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Example 2.3
5  Learning curves
6  Machine Learning for Engineering Problem Solving
7  @author: Austin Downey
8  """
9
10 import IPython as IP
11 IP.get_ipython().run_line_magic('reset', '-sf')
12
13 import numpy as np
14 import matplotlib.pyplot as plt
15 import sklearn as sk
16 from sklearn import linear_model
17 from sklearn import pipeline
18
19 plt.close('all')
20
21 """ build the data sets
22 np.random.seed(2) # 2 and 6 are pretty good
23 m = 100
24 X = 6 * np.random.rand(m,1) - 3
25 Y = 0.5 * X**2 + X + 2 + np.random.randn(m,1)
26 X_model = np.linspace(-3,3)
27
28 # plot the data
29 plt.figure()
30 plt.grid(True)
31 plt.scatter(X,Y)
32 plt.xlabel('x')
33 plt.ylabel('y')
34
35 """ generate learning curves for a linear model
36
37 # build the linear model in SK learn
38 model = sk.linear_model.LinearRegression()
39
40 # split the data into training and validation data sets
41 # Split arrays or matrices into random train and test subsets
42 X_train, X_val, y_train, y_val = sk.model_selection.train_test_split(X, Y, test_size=0.2)
43
44 train_errors, val_errors = [], []
45 for i in range(1, len(X_train)):
46     model.fit(X_train[:i], y_train[:i])
47     y_train_predict = model.predict(X_train[:i])
48     y_val_predict = model.predict(X_val)
49
50     # compute the error for the trained model
51     mse_train = sk.metrics.mean_squared_error(y_train[:i],y_train_predict)
52     train_errors.append(mse_train)
53
54     # compute the error for the validation model
55     mse_val = sk.metrics.mean_squared_error(y_val,y_val_predict)
56     val_errors.append(mse_val)
57
58 # predict model
59 y_model = model.predict(np.expand_dims(X_model,axis=1))
60
61 plt.figure('test model')
62 plt.scatter(X,Y,s=2, label='data')
63 plt.scatter(X_train[:i],y_train[:i], label='data in training set')
64 plt.scatter(X_val,y_val, marker='s', label='validation data')
65 plt.plot(X_model,y_model,'r--',label='model')
66 plt.xlabel('x')
67 plt.ylabel('y')

```

```

68     plt.legend(loc=2)
69     plt.grid(True)
70     plt.savefig('test_plots/linear_model_'+str(i))
71     plt.close('test model')
72
73     plt.figure()
74     plt.grid(True)
75     plt.plot(train_errors, "--",label="train")
76     plt.plot(val_errors, ":", label="val")
77     plt.xlabel('number of data points in training set')
78     plt.ylabel('mean squared error')
79     plt.legend(framealpha=1)
80     plt.ylim(0,6)
81     ### generate learning curves for a polynomial model
82
83     model = sk.pipeline.Pipeline((
84         ("poly_features", sk.preprocessing.PolynomialFeatures(degree=20, include_bias=False)),
85         ("lin_reg", sk.linear_model.LinearRegression()),
86     ))
87
88     # split the data into training and validation data sets
89     # Split arrays or matrices into random train and test subsets
90     X_train, X_val, y_train, y_val = sk.model_selection.train_test_split(X, Y, test_size=0.2)
91
92     train_errors = []
93     val_errors = []
94     for i in range(1, len(X_train)):
95         model.fit(X_train[:i], y_train[:i])
96         y_train_predict = model.predict(X_train[:i])
97         y_val_predict = model.predict(X_val)
98
99         # compute the error for the trained model
100        mse_train = sk.metrics.mean_squared_error(y_train[:i],y_train_predict)
101        train_errors.append(mse_train)
102
103        # compute the error for the validation model
104        mse_val = sk.metrics.mean_squared_error(y_val,y_val_predict)
105        val_errors.append(mse_val)
106
107        plt.figure('test model')
108        plt.scatter(X,Y,s=2, label='data')
109        plt.scatter(X_train[:i],y_train[:i], label='data in training set')
110        plt.scatter(X_val,y_val, marker='s', label='validation data')
111        y_model = model.predict(np.expand_dims(X_val,axis=1))
112        plt.plot(X_val,y_model,'r--',label='model')
113        plt.xlabel('x')
114        plt.ylabel('y')
115        plt.legend(loc=2)
116        plt.grid(True)
117        plt.savefig('test_plots/polynominal_model_'+str(i))
118        plt.close('test model')
119
120    plt.figure()
121    plt.grid(True)
122    plt.plot(train_errors, "--",label="train")
123    plt.plot(val_errors, ":", label="val")
124    plt.xlabel('number of data points in training set')
125    plt.ylabel('mean squared error')
126    plt.legend(framealpha=1)
127    plt.ylim(0,6)
128
129
130
131

```

Example 2.4

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Example 2.4
5  Ridge Regression
6  Machine Learning for Engineering Problem Solving
7  @author: Austin Downey
8  """
9
10 import IPython as IP
11 IP.get_ipython().run_line_magic('reset', '-sf')
12
13 import numpy as np
14 import matplotlib.pyplot as plt
15 import sklearn as sk
16 from sklearn import linear_model
17 from sklearn import pipeline
18
19 plt.close('all')
20
21 """ build the data sets
22 m = 20
23 X = 6 * np.random.rand(m, 1) - 3
24 Y = 0.5 * X**2 + X + 2 + np.random.randn(m, 1)
25
26 X_model = np.linspace(-3,3,num=1000)
27 X_model = np.expand_dims(X_model,axis=1)
28
29
30 """ Perform Ridge Regression
31
32 # plot the data
33 plt.figure()
34 plt.grid(True)
35 plt.scatter(X,Y,color='gray')
36 plt.xlabel('x')
37 plt.ylabel('y')
38
39 # build and plot a linear model
40 model_linear = sk.linear_model.Ridge(alpha=100, solver="cholesky")
41 model_linear.fit(X, Y)
42 y_model_linear = model_linear.predict(X_model)
43 plt.plot(X_model,y_model_linear,'-',label='linear model')
44
45 # build and plot a polynomial model
46 model_poly = sk.pipeline.make_pipeline(sk.preprocessing.PolynomialFeatures(10),
47                                     sk.linear_model.Ridge(alpha=100, solver="cholesky"))
48 model_poly.fit(X, Y)
49 y_model_poly = model_poly.predict(X_model)
50 plt.plot(X_model,y_model_poly,'-',label='polynomial model')
51
52 plt.legend()
53
54
55
56
57
58

```

Example 2.5

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Example 2.5
5  Early Stopping
6  Machine Learning for Engineering Problem Solving
7  @author: Austin Downey
8  """
9
10 import IPython as IP
11 IP.get_ipython().run_line_magic('reset', '-sf')
12
13 import numpy as np
14 import matplotlib.pyplot as plt
15 import sklearn as sk
16
17
18 plt.close('all')
19
20 """ build the data sets
21
22 # use 6 to help give a smooth curve that makes the case for early stopping
23 np.random.seed(6)
24
25 m = 20
26 X = 6 * np.random.rand(m, 1) - 3
27 Y = 0.5 * X**2 + X + 2 + np.random.randn(m, 1)
28
29 # plot the data
30 plt.figure()
31 plt.grid(True)
32 plt.scatter(X,Y,color='gray')
33 plt.xlabel('x')
34 plt.ylabel('y')
35
36 X_model = np.linspace(-3,3,num=1000)
37 X_model = np.expand_dims(X_model,axis=1)
38
39 X_train, X_val, y_train, y_val = sk.model_selection.train_test_split(X, Y, test_size=0.2)
40
41 """ perform early stopping
42
43
44 # prepare the data
45 poly_scaler = sk.pipeline.Pipeline([("poly_features", sk.preprocessing.PolynomialFeatures(
46     degree=90, include_bias=False)), ("std_scaler", sk.preprocessing.StandardScaler())])
47 X_train_poly_scaled = poly_scaler.fit_transform(X_train)
48 X_val_poly_scaled = poly_scaler.fit_transform(X_val)
49
50 # set up the model, not that by setting max_iter=1 it will only train one epoch
51 model = sk.linear_model.SGDRegressor(max_iter=1, tol=0, learning_rate="constant"
52     ,eta0=0.0005,penalty=None,warm_start=True)
53
54
55 # Train the model in a loop to build the data set to investigate the benefit of early
56 # stopping
57 val_errors = []
58 train_errors = []
59 for epoch in range(1000):
60     model.fit(X_train_poly_scaled, y_train.ravel()) # continues where it left off
61     y_val_predict = model.predict(X_val_poly_scaled) # Predict the target values
62     y_train_predict = model.predict(X_train_poly_scaled) # Predict the target values
63     val_error = sk.metrics.mean_squared_error(y_val, y_val_predict) # Calculate error
64     train_error = sk.metrics.mean_squared_error(y_train.ravel(), y_train_predict) #
65     Calculate error
66     val_errors.append(val_error)

```

```
65     train_errors.append(train_error)
66
67     # plot the early learning curves, you may have to plot this a few times to get
68     # a set of curves that shows strong results
69     plt.figure()
70     plt.grid(True)
71     plt.plot(val_errors, label='validation data')
72     plt.plot(train_errors, '--', label='training data')
73     plt.ylabel('RMSE')
74     plt.xlabel('epoch')
75     plt.legend()
76
77
78
79
80
81
```

3 Classification

In Chapter 1, we examined various tasks that machine learning excels at, such as regression (predicting values) and classification (identifying classes). Having explored linear and polynomial regression in Chapter 2, we now turn our attention to classification in this section.

3.1 Binary Classifier

A binary classifier is a model that assigns inputs to one of two distinct classes. Linear binary classification is a simple but powerful method used to separate data into two distinct classes by drawing a straight decision boundary—such as a line (in 2D) or hyperplane (in higher dimensions)—between them. It works by computing a weighted sum of the input features and applying a threshold to determine the output class. Mathematically, it learns a weight vector θ and bias term b such that the decision rule can be written as:

$$\hat{y} = \begin{cases} 1 & \text{if } \theta^T \mathbf{x} + b \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (3.3)$$

where \mathbf{x} is the input feature vector. During training, the weights are adjusted to minimize classification error, often using gradient descent on a suitable loss functions (e.g., hinge loss or logistic loss).

Despite its simplicity, the linear classifier is often surprisingly effective, especially when the data is linearly separable or nearly so. An example of such a case is shown in figure 3.1 where the three models (H_1 , H_2 , and H_3) are linear models defined by the hyperparameters θ that each correctly classify the data. The challenge then becomes how does one obtain the hyperparameters θ . To start, we will show that the hyperparameters θ can be obtained through gradient descent.

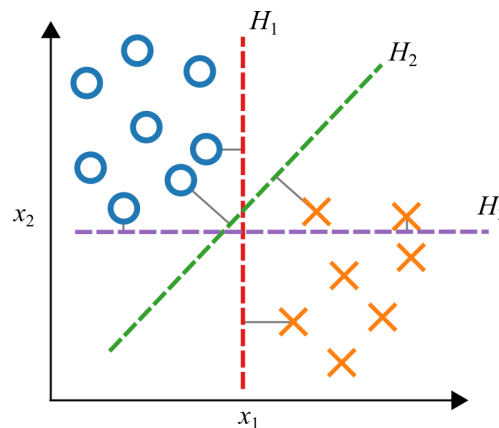


Figure 3.1: Linear classifier that can be implemented with gradient descent.

Review 3.1 MNIST (Modified National Institute of Standards and Technology) Dataset

The MNIST (Modified National Institute of Standards and Technology) dataset is a collection of 70,000 small images of handwritten digits (figure 3.2). This dataset is known as “modified” because it combines two earlier sets of images: one created by U.S. high school students and another by US Census Bureau employees. First published in 1998, the dataset is particularly suitable for machine learning tasks as each image is labeled with the digit it represents. Each digit has been normalized and centered in a 28x28 pixel frame and anti-aliased to introduce grayscale levels. Each instance (i.e. picture) has 784 features that range from 0 to 256 and each represent an 8-bit gray scale pixel as shown in figure 3.3.

Typically, the dataset is split into two parts: the first 60,000 images are used for training, and the remaining 10,000 serve as validation data. It is standard practice to train and test models on the first 60,000 images, reserving the last 10,000 for final validation at the end of the project when the classifier is ready for deployment. Scikit-Learn includes a helper function that facilitates the downloading of the MNIST dataset.

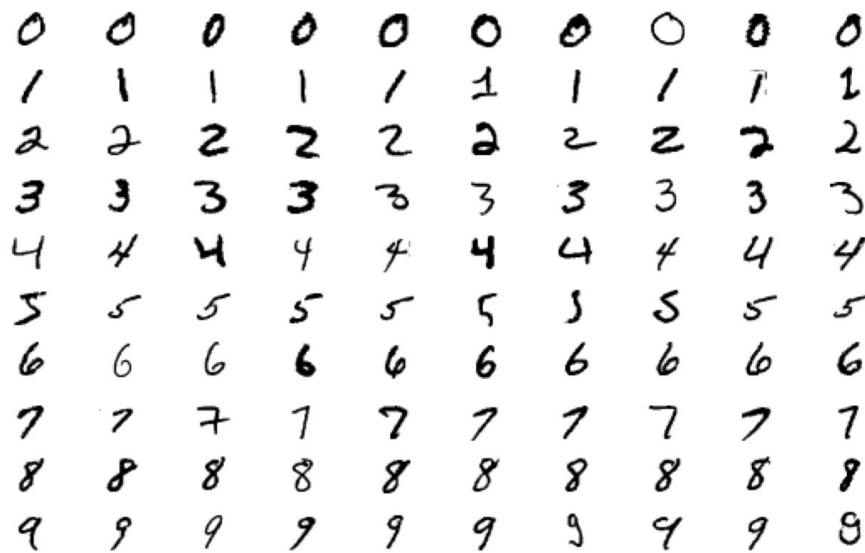
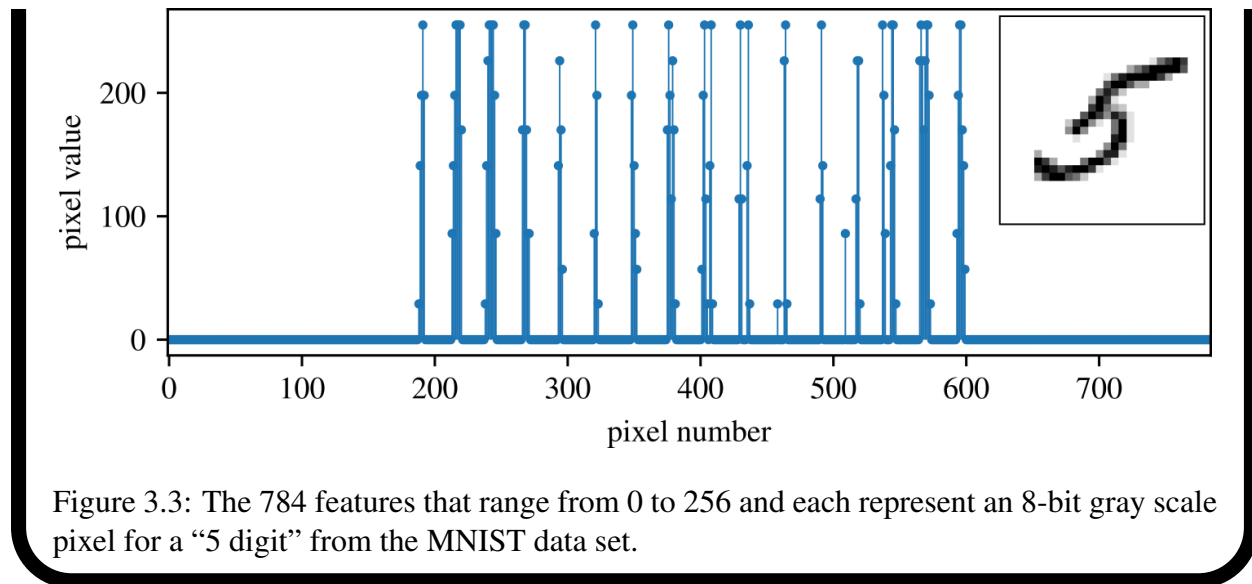


Figure 3.2: A collection of the 10 digits (0-9) that make up the MNIST data set.



To begin, we simplify our task by focusing on identifying all the “5”s in the MNIST dataset. This task involves using a binary classifier that checks whether each digit is a 5. The number 5 is notoriously difficult to classify correctly within the MNIST dataset, making it a sensible starting point for binary classification.

Example 3.1 MNIST Dataset

This example demonstrates how to use `sklearn.datasets.fetch_openml` to load the MNIST dataset. It also includes visualization steps to display sample handwritten digits.

3.1.1 Regularized Linear Classifier

A practical algorithm for building this “5-detector” is the Regularized Linear Classifier solved with Stochastic Gradient Descent (SGD). The Regularized Linear Classifier offers an efficient and straightforward approach for discriminative learning of linear classifiers under convex loss functions, scaling well with large datasets and maintaining simplicity in implementation as it processes one training example at a time. The main advantages of the Regularized Linear Classifier include:

- Efficiency.
- Ease of implementation, providing many opportunities for optimization.

However, the Stochastic Gradient Descent classifier also presents certain disadvantages:

- It requires careful tuning of several hyperparameters, including the regularization parameter and the number of iterations.
- It is sensitive to feature scaling.

NOTE

The Regularized Linear Classifier solved with Stochastic Gradient Descent is commonly referred to simply called a Stochastic Gradient Descent classifier, despite Stochastic Gradient Descent only being the method used to train the linear model. This is exemplified by the fact that scikit-learn uses the name `SGDClassifier` for their model that implements a “Regularized linear models with stochastic gradient descent (SGD) learning”.

Given a labeled training set $(\mathbf{x}^{(i)}, y^{(i)})_{i=1}^m$ with $y^{(i)} \in \{1, 0\}$ (1 marks the digit “5”), the regularised linear “5-detector” learns a weight vector θ (bias in the first entry θ_1) by minimising

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \max(0, 1 - y^{(i)} \theta^T \mathbf{x}^{(i)}) + \frac{\lambda}{2} \|\theta_{1:}\|_2^2, \quad (3.4)$$

where the first term is the hinge loss enforcing a unit margin around the decision boundary. The second term applies ℓ_2 -regularisation of strength $\lambda > 0$ to all weights except the bias (θ_1). Minimising (3.4) with stochastic-gradient descent yields a sparse, margin-maximising classifier that generalises well to unseen handwritten digits.

Example 3.2 Regularized Linear Classifier

This example trains a regularized linear classifier using `SGDClassifier` to distinguish the digit “5” from other digits in the MNIST dataset. It highlights how regularization helps improve generalization performance by penalizing large weights, and demonstrates basic prediction on a selected digit using a learned linear decision boundary.

3.2 Performance Measures for Binary Classification

Assessing the performance of a classifier can be more intricate than evaluating a regressor. Nonetheless, there are numerous performance metrics available that facilitate the comprehensive evaluation of a classifier’s effectiveness.

3.2.1 Confusion Matrix

First, a brief explanation of Type I and Type II errors. These terms are not exclusive to classification problems in machine learning; they are also crucial in statistical hypothesis testing:

Type I Error: False positive (incorrectly rejecting a true null hypothesis)

Type II Error: False negative (incorrectly failing to reject a false null hypothesis)

These four cases from statistical testing can be combined into a single matrix known as the confusion matrix. In the context of our 5-detector, consider the simplified confusion matrix shown in figure 3.4.

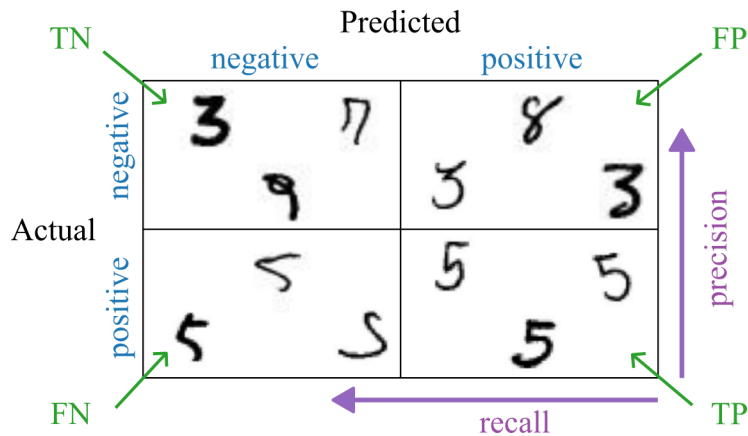


Figure 3.4: Confusion matrix binary.

A effective approach to evaluating classifier performance is to examine the confusion matrix. The essence is to tally the instances where class A is incorrectly classified as class B. A well-performing classifier will exhibit high values along the diagonal and lower values elsewhere in the matrix. Each row represents the actual class, while each column represents the predicted class. For instance, to determine how often the classifier misidentified images of 5s as 3s, inspect the intersection of the 5th row and 3rd column in figure 3.5.

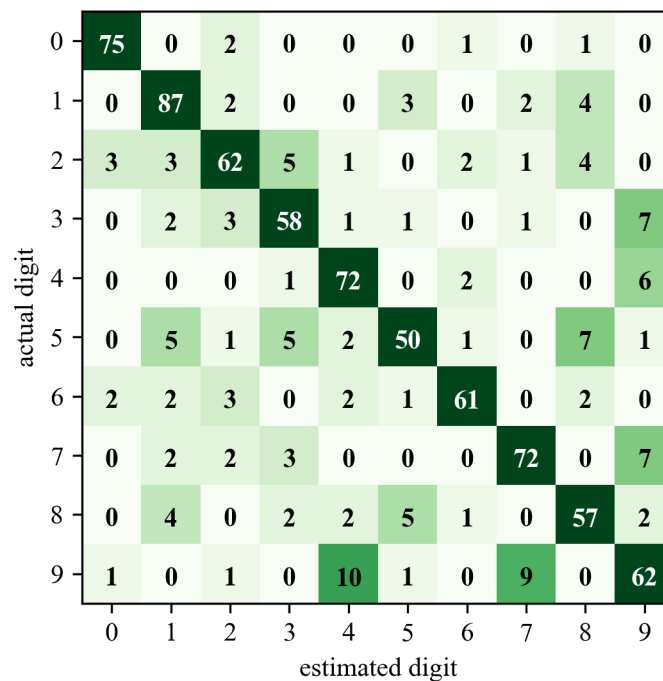


Figure 3.5: Confusion matrix for the first 1,000 data points from the MNIST dataset.

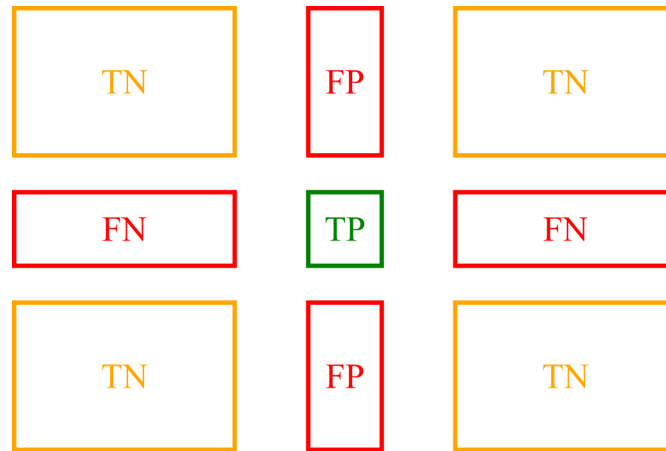


Figure 3.6: Confusion matrix breakout.

Each row in the confusion matrix corresponds to an actual class, while each column corresponds to a predicted class. This is diagrammed in figure 3.6. The first row represents non-5 images (the negative class), and the second row represents images classified as 5s (the positive class). The first column of the first row indicates true negatives (TN), while the second column of the first row shows false positives (FP). The second row denotes images of 5s (the positive class): the first column represents false negatives (FN), and the second column displays true positives (TP). A perfect classifier would only have nonzero values along the main diagonal (from top left to bottom right).

Example 3.3 Binary Confusion Matrix

This example evaluates a binary classifier trained to detect the digit “5” using a confusion matrix. The matrix summarizes the number of true positives, true negatives, false positives, and false negatives. It also visualizes specific misclassifications to help interpret where the model fails, such as incorrectly labeling a “5” as not a “5” (false negative) or misidentifying another digit as a “5” (false positive).

3.2.2 Accuracy, Precision, and Recall

The confusion matrix provides detailed information, but sometimes a more concise metric is preferred. Before looking into accuracy, let’s define it as

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total number of Predictions}} \quad (3.5)$$

or

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (3.6)$$

While accuracy offers valuable insights, it’s essential to consider the number of falsely classified positive and negative predictions, especially within the context of the problem being addressed. For instance, a 99% accuracy rate might seem satisfactory for predicting credit card fraud. However, what if a false negative indicates a severe virus or cancer? This is why we need to look deeper into the accuracy formula.

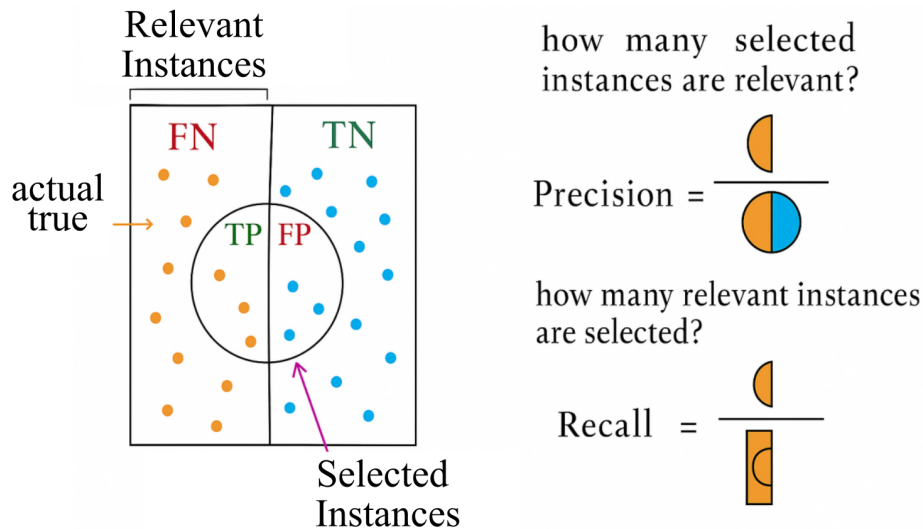


Figure 3.7: Visual representation of relevant instances, showing their relation to True Positives (TP) and False Positives (FP).

This is where precision and recall come into play. For this, we will need a definitions of instances

$$\text{Selected instances} = TP + FP \quad (3.7)$$

and

$$\text{Relevant instances} = TP + FN \quad (3.8)$$

which are visualized in figure 3.7. With these definitions in mind in mind, let's define precision as

$$\text{Precision} = \frac{\text{selected relevant instances}}{\text{selected instances}}, \quad (3.9)$$

or,

$$\text{Precision} = \frac{TP}{TP + FP}. \quad (3.10)$$

One way to achieve perfect precision is to make only one positive prediction and ensure it's correct (precision = 1/1 = 100%). However, this approach would be impractical as it would ignore most positive instances. Recall complements precision by considering all relevant results, both true positives and false negatives. Recall, also known as sensitivity or true positive rate (TPR), is the ratio of positive instances correctly detected by the classifier. We can define recall as

$$\text{Recall} = \frac{\text{selected relevant instances}}{\text{relevant instances}}, \quad (3.11)$$

or,

$$\text{Recall} = \frac{TP}{TP + FN}. \quad (3.12)$$

While the definition of accuracy, precision, and recall may not be immediately intuitive, a visualizations such as that shown in figure 3.7 can help clarify them.

To balance precision and recall, we often use the F1 score, which is the harmonic mean of precision and recall. The F1 score favors classifiers with similar precision and recall, computed as

$$F_1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} = \frac{TP}{TP + \frac{FN+FP}{2}}. \quad (3.13)$$

However, optimizing both precision and recall simultaneously is challenging due to the precision/recall tradeoff.

Figure 3.8 illustrates the tradeoffs when adjusting the decision threshold in a Regularized Linear Classifier. By adjusting the classification threshold, we can control the balance between precision and recall. Lowering the threshold increases recall but reduces precision, and vice versa.

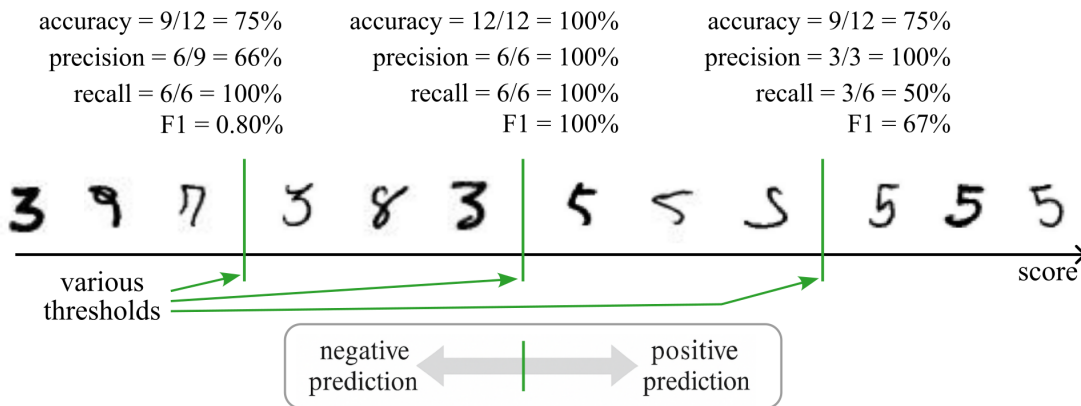


Figure 3.8: Visualization changing thresholds on the classification of digits in a “5-detector” and the resulting change in precision and recall.

To determine the optimal threshold, it’s helpful to plot precision and recall against the threshold values. By default, the classifier in scikit-learn uses a threshold of zero, but this can be adjusted as needed. As illustrated in Figure 3.9, it’s relatively straightforward to develop a classifier with nearly any desired precision: simply adjust the threshold to a sufficiently high value. However, it’s important to bear in mind that a high-precision classifier may not be very practical if its recall is too low.

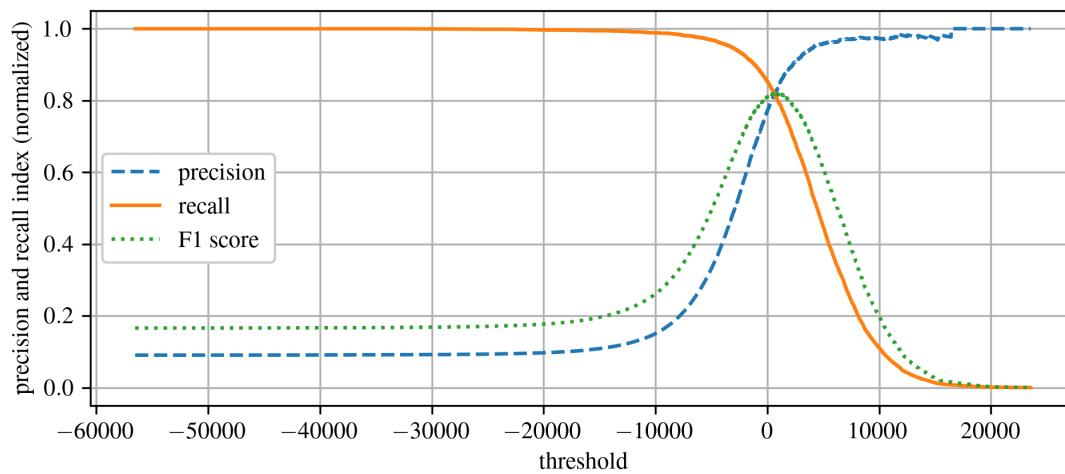


Figure 3.9: Visualization of the precision, recall, and F1 score as a function of changing a threshold value.

Figure 3.10 reports a precision vs recall curve that is obtained by changing the threshold of the trained model. At the leftmost part of the precision-recall curve (low recall), precision is highest because the model makes very few positive predictions, and those are likely true positives. As recall increases, the model starts predicting more positives, including more false positives, which reduces precision. The shape of the curve is also important; a steep drop-off indicates that the addition of false positives happens rapidly with increasing recall, whereas a more gradual decline suggests a better balance between precision and recall.

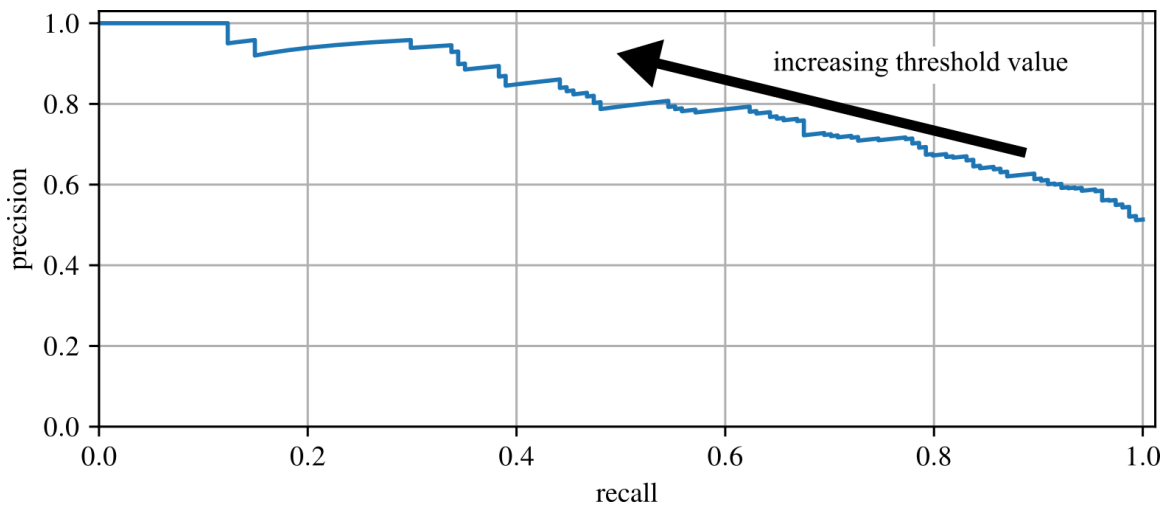


Figure 3.10: The relationship between precision vs recall for a changing threshold value.

By analyzing the curve, one can choose a threshold that balances precision and recall according to the specific needs of the application. For instance, in medical diagnosis, high recall is crucial to ensure no cases are missed, even if precision is lower. Additionally, the area under the precision-recall curve (AUC-PR) can be used to compare models. A higher area indicates a model with better performance across all thresholds.

Example 3.4 Accuracy, Precision, and Recall

This example computes and compares three key classification metrics (accuracy, precision, and recall) for a linear classifier trained to detect the digit “5” in the MNIST dataset. It shows both manual and `sklearn`-based calculations, introduces the F1 score, and visualizes how precision and recall vary with the classification threshold using a precision-recall curve.

3.3 k -fold Cross-validation

As you've likely observed, training your SGD classifier doesn't always yield consistent results in terms of model performance; hence the "stochastic" in stochastic gradient descent. One approach to obtaining a solution closer to the global minimum is to allow the gradient descent algorithm to iterate over more epochs. However, this poses a challenge as it's highly time-consuming and may lead to overfitting the model. Moreover, the variability in model performance makes it difficult to gauge the model's effectiveness.

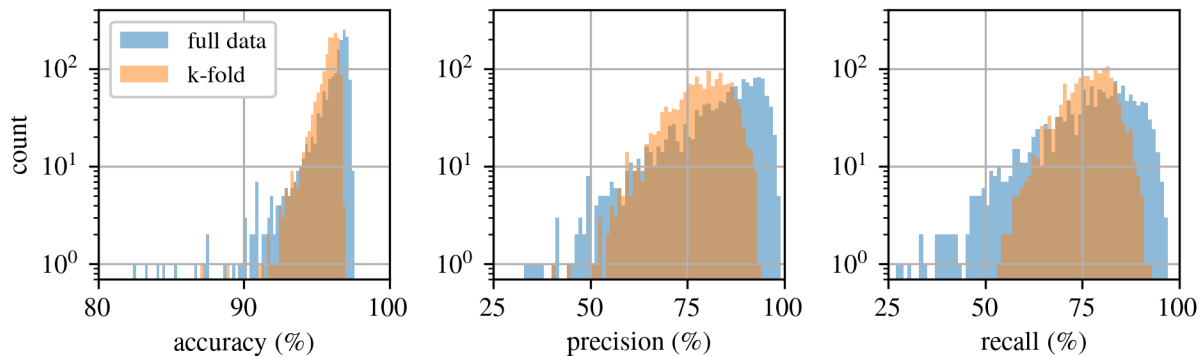


Figure 3.11: Performance metrics for 500 5-detectors trained either with full data or through k -fold cross-validation with 3 folds.

k -fold cross-validation offers a more rigorous method for evaluating your algorithm's performance. For instance, consider Figure 3.11, which depicts the performance metrics for 500 classification models trained as 5-detectors. Noticeably, the results obtained using k -fold cross-validation exhibit significantly less variability compared to training the classifier on the full dataset. In k -fold cross-validation, the training set is divided into smaller subsets for training and validation. The model is trained on these subsets and evaluated on the validation sets. Figure 3.12 provides a graphical representation of this technique.

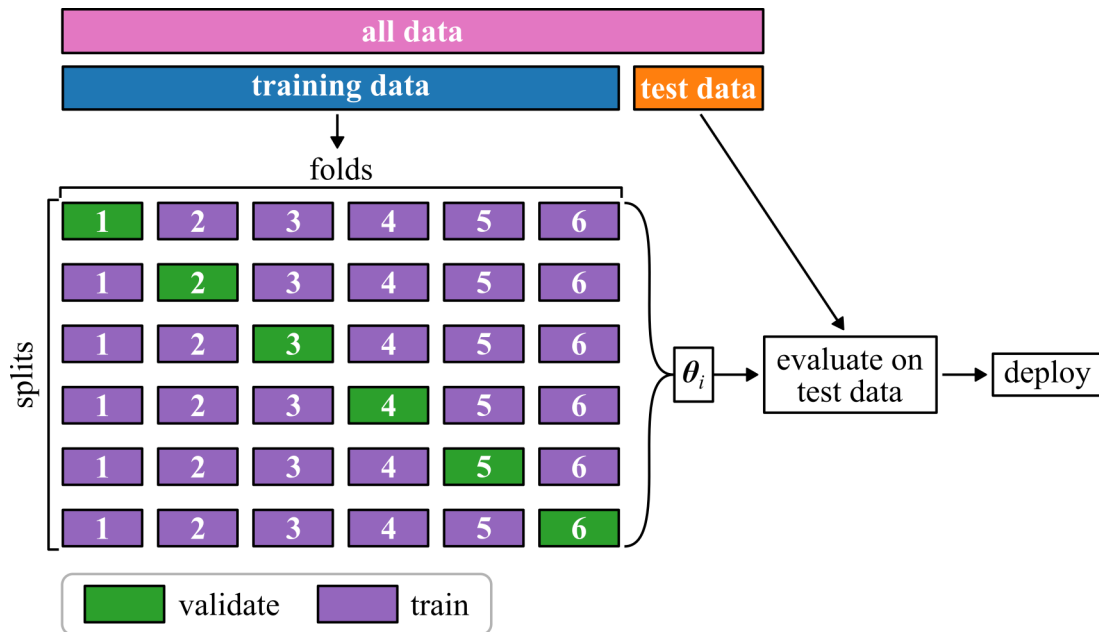


Figure 3.12: Illustration of how k -fold cross-validation operates.

Of course, Scikit-Learn provides a cross-validation feature, known as `sk.model_selection.cross_val_score`. This function conducts k -fold cross-validation by randomly partitioning the training set into distinct folds. Each fold is used for validation while the rest are used for training the SGD classifier. The reported performance metrics are averaged across all folds. It's important to note that by partitioning the data into folds, the number of samples available for learning is limited, which may result in peculiar classifiers with unusually high or low error rates. Despite being more computationally expensive, this approach effectively utilizes all the data for training and validation.

Example 3.5 k -fold Cross-Validation

This example trains a binary classifier to detect the digit “5” from the MNIST dataset using Stochastic Gradient Descent (SGD). It first demonstrates how model performance can vary when trained multiple times on the same dataset. Then, by applying k -fold cross-validation using `sk.model_selection.cross_val_predict`, it shows how splitting the data into k subsets helps reduce this variation and provides a more reliable estimate of accuracy, precision, and recall.

3.4 Multiclass Classification

Previously, we explored binary classifiers, which distinguish between two classes. Now, we'll explore the use of multiclass classifiers (also known as multinomial classifiers), which can discern between more than two classes. The simplest approach to multiclass classification involves multiple iterations of a binary classifier. There are two primary strategies for achieving this, both of which are considered in the context of the MNIST database challenge:

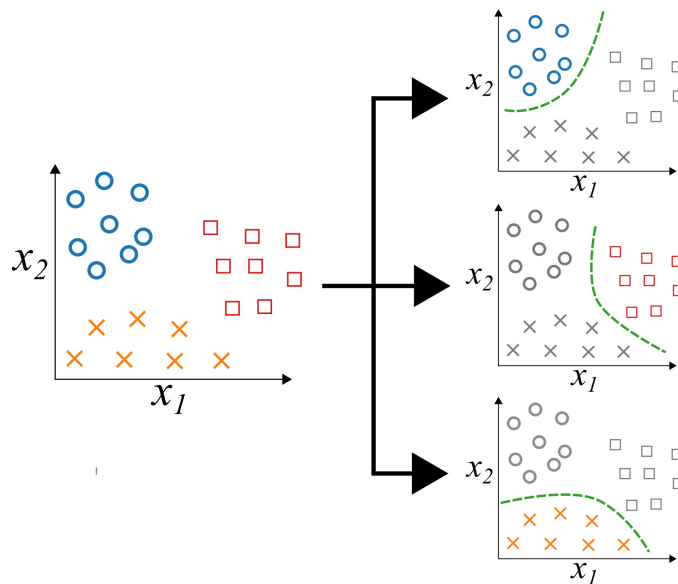


Figure 3.13: One-versus-Rest (OvR) classifier for a multi-class dataset.

- **One-versus-Rest (OvR):** This strategy involves creating a system capable of classifying digit images into 10 classes (from 0 to 9) by training 10 binary classifiers, one for each digit (e.g., a 0-detector, a 1-detector, etc.). When classifying an image, the decision score from each classifier is obtained, and the class with the highest score is selected. While sometimes referred to as One-Versus-All (OvA), this term is slightly misleading as it does not involve comparisons against the same class.

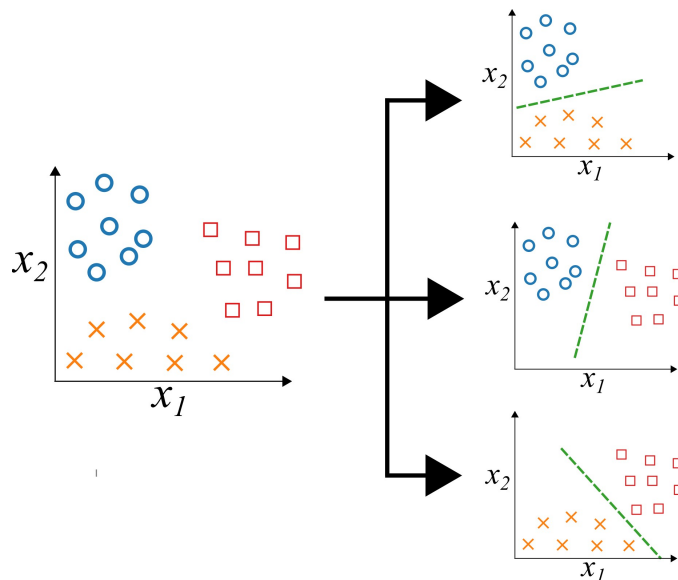


Figure 3.14: One-versus-one (OvO) classifier for a multi-class dataset.

- **One-versus-one (OvO):** This approach trains a binary classifier for every pair of digits, such as one for distinguishing between 0s and 1s, another for distinguishing between 0s and

5s, and so forth. If there are N classes, $N \times (N - 1)/2$ classifiers need to be trained. For the MNIST problem, this translates to training 45 binary classifiers. During classification, the image is evaluated against all 45 classifiers, and the class with the most victories is selected. The primary advantage of OvO is that each classifier only needs to be trained on the relevant portion of the training set for the two classes it distinguishes, which is beneficial for algorithms that scale poorly.

Certain algorithms, such as Support Vector Machine classifiers, suffer from scalability issues with larger training sets. Algorithms that do not scale well with large datasets often resort to the OvO approach because training many small classifiers is computationally cheaper than fitting a handful of models on the full data. For most binary classifiers the OvR strategy is the norm. Some methods (including Random Forests and naïve Bayes) support multiclass problems natively and therefore require neither reduction strategy.

Example 3.6 Multiclass Regularized Linear Classifier

This example extends the use of `SGDClassifier` to multiclass classification on the MNIST dataset. It compares three approaches: One-vs-Rest (OvR), One-vs-One (OvO), and Scikit-Learn's default automatic strategy. Each method trains multiple binary classifiers to distinguish between the 10 digit classes.

3.5 Performance Measures for Multiclass Classification

Assessing a multiclass classifier often poses more challenges compared to evaluating a binary classifier. Nonetheless, many of the principles and methodologies utilized in evaluating binary classifiers can be seamlessly applied to multiclass classifiers as well.

3.5.1 Confusion Matrix

First, examine the confusion matrix. To do this, generate predictions using the `cross_val_predict()` function, followed by calling the `confusion_matrix()` function, similar to your previous steps. The confusion matrix in figure 3.15 appears satisfactory, with most images located on the main diagonal, indicating correct classification. However, the 5s appear slightly darker than other digits, suggesting fewer instances of 5s in the dataset or inferior performance of the classifier on 5s. Further investigation confirms both scenarios.

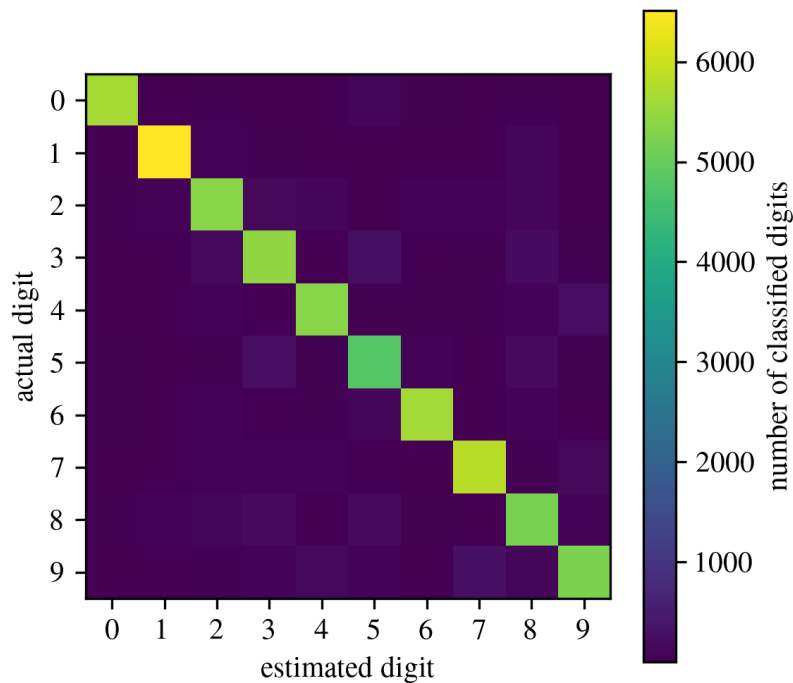


Figure 3.15: Confusion matrix for the MNIST data set solved using a one-vs-one classifier with Stochastic Gradient Descent and a k -fold of 3.

To turn figure 3.15 into a figure that highlights the mistakes, first normalize each entry of the confusion matrix by the number of images in its true class so you compare error rates rather than raw counts that would overweight frequent classes. Then set the diagonal cells to NaN so only the off-diagonal errors remain visible, and plot the resulting matrix. The resulting figure 3.16 more clearly highlights the errors in the confusion matrix.

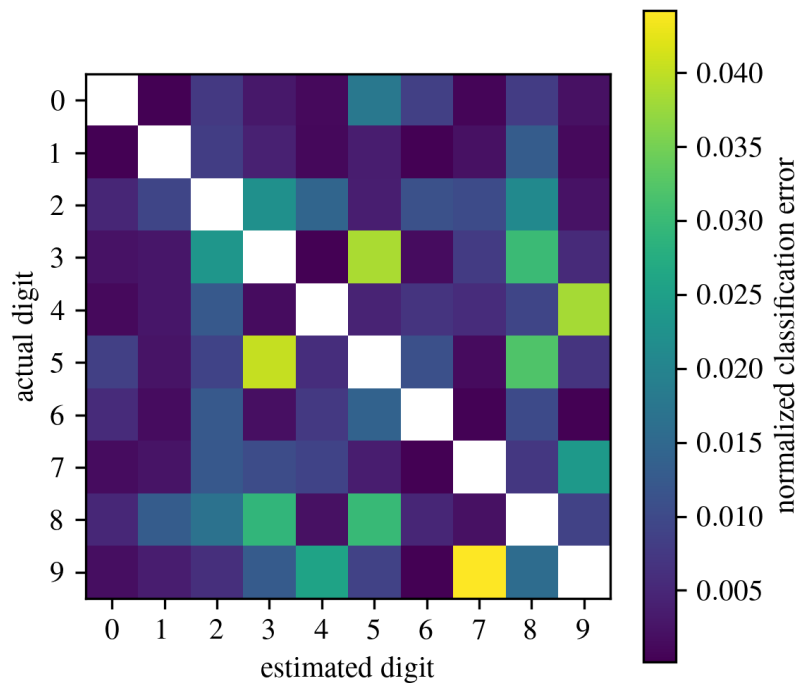


Figure 3.16: Normalized confusion matrix (similar to figure 3.15) with the diagonal removed to emphasize errors.

The refined plot facilitates a clear understanding of the classifier's error patterns. Rows represent actual classes, while columns depict predicted classes. Columns corresponding to classes 8 and 9 exhibit brightness, indicating numerous misclassifications as 8s or 9s. Likewise, rows for classes 8 and 9 also appear bright, signifying frequent confusion between 8s, 9s, and other digits. Conversely, some rows, like row 1, display darkness, indicating accurate classification of most 1s (with few exceptions confused with 8s). Notably, errors are asymmetric; for instance, more 5s are misclassified as 8s than vice versa.

Analyzing the confusion matrix provides valuable insights for classifier enhancement. In this case, efforts should concentrate on improving the classification of 8s and 9s, along with addressing the specific 3/5 confusion. Potential strategies include gathering more training data for these digits, devising new features (e.g., counting closed loops), or preprocessing images to emphasize distinguishing patterns (e.g., using Scikit-Image, Pillow, or OpenCV).

3.5.2 Analyzing Individual Errors

Analyzing individual errors provides valuable insights into the classifier's behavior and reasons for its failures.

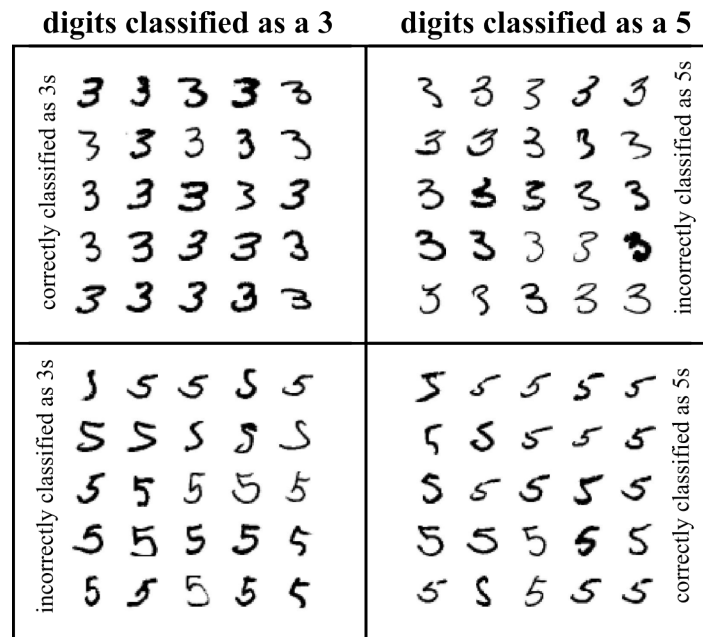


Figure 3.17: Confusion matrix showing digits classified as 3s and 5s for the same classifier used to obtain the results shown in figure 3.15.

The left half of figure 3.17 contains the blocks for images the classifier labelled as 3, while the right half shows those it labelled as 5. A few misclassified digits are so poorly written that even a human might hesitate, yet many errors look quite clear. Remember that a regularised linear classifier assigns a single weight to every pixel for each class and then sums the weighted pixel intensities to score each class.

Because just a handful of pixels separate a 3 from a 5, a linear classifier often confuses the two. Their main distinction is the small stroke that joins the top bar to the bottom curve; even a slight shift or rotation of this junction can flip the prediction. The linear model is therefore very sensitive to translations and rotations. Pre-processing the images to center the digits and correct their orientation should lessen the 3-5 confusion and improve performance across the board. 5

A non-linear model can learn more complex patterns by applying transformations (such as kernel mappings), effectively giving different relative importance to various parts of the image, and would therefore be better equipped to capture more intricate pixel patterns and more effectively separate these digits.

Example 3.7 Multiclass Confusion Matrix

This example uses a One-vs-One classifier trained with Stochastic Gradient Descent to evaluate multiclass classification performance on the MNIST dataset. It constructs a confusion matrix using 3-fold cross-validation and visualizes both the raw classification counts and the normalized error rates, helping identify which digits are most commonly misclassified.

3.6 Examples

Example 3.1

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Example 3.1 Load the MNIST data set
5  Created for EMCH 504 at USC
6  @author: Austin Downey
7  """
8
9  import IPython as IP
10 IP.get_ipython().run_line_magic('reset', '-sf')
11
12 import numpy as np
13 import scipy as sp
14 import matplotlib as mpl
15 import matplotlib.pyplot as plt
16 import sklearn as sk
17 from sklearn import linear_model
18 from sklearn import datasets
19
20 plt.close('all')
21
22
23 #%% Load your data
24
25 # this fetches "a" MNIST dataset from openml and loads it into your environment
26 # as a Bunch, a Dictionary-like object that exposes its keys as attributes.
27 mnist = sk.datasets.fetch_openml('mnist_784',as_frame=False,parser='auto')
28
29 # calling the DESCR key will return a description of the dataset
30 print(mnist['DESCR'])
31
32 # calling the data key will return an array with one row per instance and one
33 # column per feature where each features is a pixel, as defined in the key feature_names
34 X = mnist['data']
35
36
37 # calling the target key will return an array with the labels
38 Y = np.asarray(mnist['target'],dtype=int)
39
40 # Each image is 784 features or 28x28 pixels, however, the features must be reshaped
41 # into a 29x29 grid to make them into a digit, where the values represents one
42 # the intensity of one pixel, from 0 (white) to 255 (black).
43
44 digit_id = 35 # An OK 5
45 # digit_id = 0 # An odd 5
46 # digit_id = 100 # A bad 5
47
48
49 test_digit = X[digit_id,:]
50 digit_resaped = np.reshape(test_digit,(28,28))
51
52 # plot an image of the random pixel you picked above.
53 plt.figure()
54 plt.imshow(digit_resaped,cmap = mpl.cm.binary,interpolation="nearest")
55 plt.title('A '+str(Y[digit_id])+' digit from the MNIST dataset')
56 plt.xlabel('pixel column number')
57 plt.ylabel('pixel row number')
58
59

```

Example 3.2

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Example 3.3 Binary confusion matrix for the MNIST dataset
5  Created for EMCH 504 at USC
6  @author: Austin Downey
7  """
8
9  import IPython as IP
10 IP.get_ipython().run_line_magic('reset', '-sf')
11
12 import numpy as np
13 import scipy as sp
14 import matplotlib as mpl
15 import matplotlib.pyplot as plt
16 import sklearn as sk
17 from sklearn import linear_model
18 from sklearn import pipeline
19 from sklearn import datasets
20 from sklearn import metrics
21
22 cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
23 plt.close('all')
24
25
26 """ Load your data
27
28 # Fetch the MNIST dataset from openml
29 mnist = sk.datasets.fetch_openml('mnist_784', as_frame=False, parser='auto')
30 X = mnist['data'] # load the data
31 Y = np.asarray(mnist['target'], dtype=int) # load the target
32
33 # Split the data set up into a training and testing data set
34 X_train = X[0:60000,:]
35 X_test = X[60000:,:]
36 Y_train = Y[0:60000]
37 Y_test = Y[60000:]
38
39 """ Train a Stochastic Gradient Descent classifier
40
41 # Extract a subset for our "5-dector".
42 Y_train_5 = (Y_train == 5)
43 Y_test_5 = (Y_test == 5)
44
45 # build and train the classifier
46 sgd_clf = sk.linear_model.SGDClassifier()
47 sgd_clf.fit(X_train, Y_train_5)
48
49 # we can now test this for the "5" that we plotted earlier.
50 digit_id = 35
51 test_digit = X[digit_id,:]
52 print(sgd_clf.predict([test_digit]))
53
54 """ Build the Confusion Matrices
55
56 # Return the predictions made on each test fold
57 X_train_pred = sgd_clf.predict(X_train)
58
59
60 # build the confusion Matrix
61 print(sk.metrics.confusion_matrix(Y_train_5, X_train_pred))
62
63 # Now let's find all the False positive and false negative
64 confusion_booleans = np.vstack((Y_train_5, X_train_pred)).T
65 FN_index = np.where((confusion_booleans == [True, False]).all(axis=1))[0]
66 FP_index = np.where((confusion_booleans == [False, True]).all(axis=1))[0]
67

```

```
68 # We built a 5-detector, so:
69 # True and True is true positive (TP)
70 # False and False is True Negative (TN)
71 # True and False is false negative (FN)
72 # False and True is false positive (FP)
73
74
75 # from this, we see #0 is a false negative (FN), i.e. its is actually a 5 but
76 # the classifier said it was not. We can plot this digit below
77 digit_id = 0
78 test_digit = X[digit_id,:]
79 digit_resaped = np.reshape(test_digit,(28,28))
80 plt.figure()
81 plt.imshow(digit_resaped,cmap = mpl.cm.binary,interpolation="nearest")
82 plt .title('A "'+str(Y[digit_id])+'" digit from the MNIST dataset')
83 plt.xlabel('pixel column number')
84 plt.ylabel('pixel row number')
85
86 # Lastly, we see that #68 is classified as an false positive (FP), again, we can plot
87 this.
88 digit_id = 161
89 test_digit = X[digit_id,:]
90 digit_resaped = np.reshape(test_digit,(28,28))
91 plt.figure()
92 plt.imshow(digit_resaped,cmap = mpl.cm.binary,interpolation="nearest")
93 plt .title('A "'+str(Y[digit_id])+'" digit from the MNIST dataset')
94 plt.xlabel('pixel column number')
95 plt.ylabel('pixel row number')
96
97
```

Example 3.3

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Example 3.3 Binary confusion matrix for the MNIST dataset
5  Created for EMCH 504 at USC
6  @author: Austin Downey
7  """
8
9  import IPython as IP
10 IP.get_ipython().run_line_magic('reset', '-sf')
11
12 import numpy as np
13 import scipy as sp
14 import matplotlib as mpl
15 import matplotlib.pyplot as plt
16 import sklearn as sk
17 from sklearn import linear_model
18 from sklearn import pipeline
19 from sklearn import datasets
20 from sklearn import metrics
21
22 cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
23 plt.close('all')
24
25
26 """ Load your data
27
28 # Fetch the MNIST dataset from openml
29 mnist = sk.datasets.fetch_openml('mnist_784', as_frame=False, parser='auto')
30 X = mnist['data'] # load the data
31 Y = np.asarray(mnist['target'], dtype=int) # load the target
32
33 # Split the data set up into a training and testing data set
34 X_train = X[0:60000,:]
35 X_test = X[60000:,:]
36 Y_train = Y[0:60000]
37 Y_test = Y[60000:]
38
39 """ Train a Stochastic Gradient Descent classifier
40
41 # Extract a subset for our "5-dector".
42 Y_train_5 = (Y_train == 5)
43 Y_test_5 = (Y_test == 5)
44
45 # build and train the classifier
46 sgd_clf = sk.linear_model.SGDClassifier()
47 sgd_clf.fit(X_train, Y_train_5)
48
49 # we can now test this for the "5" that we plotted earlier.
50 digit_id = 35
51 test_digit = X[digit_id,:]
52 print(sgd_clf.predict([test_digit]))
53
54 """ Build the Confusion Matrices
55
56 # Return the predictions made on each test fold
57 X_train_pred = sgd_clf.predict(X_train)
58
59
60 # build the confusion Matrix
61 print(sk.metrics.confusion_matrix(Y_train_5, X_train_pred))
62
63 # Now let's find all the False positive and false negative
64 confusion_booleans = np.vstack((Y_train_5, X_train_pred)).T
65 FN_index = np.where((confusion_booleans == [True, False]).all(axis=1))[0]
66 FP_index = np.where((confusion_booleans == [False, True]).all(axis=1))[0]
67

```

```
68 # We built a 5-detector, so:
69 # True and True is true positive (TP)
70 # False and False is True Negative (TN)
71 # True and False is false negative (FN)
72 # False and True is false positive (FP)
73
74
75 # from this, we see #0 is a false negative (FN), i.e. its is actually a 5 but
76 # the classifier said it was not. We can plot this digit below
77 digit_id = 0
78 test_digit = X[digit_id,:]
79 digit_resaped = np.reshape(test_digit,(28,28))
80 plt.figure()
81 plt.imshow(digit_resaped,cmap = mpl.cm.binary,interpolation="nearest")
82 plt .title('A "'+str(Y[digit_id])+'" digit from the MNIST dataset')
83 plt.xlabel('pixel column number')
84 plt.ylabel('pixel row number')
85
86 # Lastly, we see that #68 is classified as an false positive (FP), again, we can plot
87 this.
88 digit_id = 161
89 test_digit = X[digit_id,:]
90 digit_resaped = np.reshape(test_digit,(28,28))
91 plt.figure()
92 plt.imshow(digit_resaped,cmap = mpl.cm.binary,interpolation="nearest")
93 plt .title('A "'+str(Y[digit_id])+'" digit from the MNIST dataset')
94 plt.xlabel('pixel column number')
95 plt.ylabel('pixel row number')
96
97
```

Example 3.4

```

1  """
2  Example 3.4 Precision and recall accuracy for the MNIST dataset
3  @author: austin_downey
4  """
5
6  import IPython as IP
7  IP.get_ipython().run_line_magic('reset', '-sf')
8
9  import numpy as np
10 import scipy as sp
11 import matplotlib.pyplot as plt
12 import sklearn as sk
13 from sklearn import linear_model
14 from sklearn import datasets
15 from sklearn import metrics
16
17 cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
18 plt.close('all')
19
20
21 %% Load your data
22
23 # Fetch the MNIST dataset from openml
24 mnist = sk.datasets.fetch_openml('mnist_784',as_frame=False,parser='auto')
25 X = np.asarray(mnist['data']) # load the data
26 Y = np.asarray(mnist['target'],dtype=int) # load the target
27
28 # Split the data set up into a training and testing data set
29 X_train = X[0:60000,:]
30 X_test = X[60000:,:]
31 Y_train = Y[0:60000]
32 Y_test = Y[60000:]
33
34 %% Train a Stochastic Gradient Descent classifier
35
36 # Extract a subset for our "5-detector".
37 Y_train_5 = (Y_train == 5)
38 Y_test_5 = (Y_test == 5)
39
40 # build and train the classifier
41 sgd_clf = sk.linear_model.SGDClassifier()
42 sgd_clf.fit(X_train, Y_train_5)
43
44 %% Build the Confusion Matrices
45
46 # Return the predictions made with the trained model
47 X_train_pred = sgd_clf.predict(X_train)
48
49 # build the confusion Matrix
50 confusion_matrix = sk.metrics.confusion_matrix(Y_train_5, X_train_pred)
51
52 TN = confusion_matrix[0,0]
53 FP = confusion_matrix[0,1]
54 FN = confusion_matrix[1,0]
55 TP = confusion_matrix[1,1]
56
57 %% Calculate Precision and Recall
58
59 # calculate the Precision and Recall values using the commands discussed in class
60 accuracy = (TP + TN)/(TP + TN + FP + FN)
61 precision = TP/(TP+FP)
62 recall = TP/(TP+FN)
63
64 # of course, SK learn has built-in functions for this.
65 accuracy_SK = sk.metrics.accuracy_score(Y_train_5, X_train_pred)
66 precision_SK = sk.metrics.precision_score(Y_train_5, X_train_pred)
67 recall_SK = sk.metrics.recall_score(Y_train_5, X_train_pred)

```

```
68
69 # compute the F1 score for the data set
70 F1 = sk.metrics.f1_score(Y_train_5, X_train_pred)
71
72 %% plot the precision and recall over the threshold domain.
73
74 # first, compute the scores for the predictions.
75 Y_scores= sgd_clf.decision_function(X_train)
76
77 #Y_scores= Y3
78 # now, for the scores and the target training set, calculate the precisions, recalls,
79 # threshold values.
80 precisions, recalls, thresholds = sk.metrics.precision_recall_curve(Y_train_5, Y_scores)
81
82 # now, compute the F1 score over the entire threshold range
83 F1s = 2/(1/precisions + 1/recalls)
84
85 # plot the Precision and Recall vs threshold.
86 plt.figure()
87 plt.plot(thresholds, precisions[:-1], "--", label="Precision")
88 plt.plot(thresholds, recalls[:-1], "-", label="Recall")
89 plt.plot(thresholds, F1s[:-1], ":", label="F1 score")
90 plt.xlabel("Threshold")
91 plt.ylabel("normalized precision\nand recall index")
92 plt.legend(loc=6,framealpha=1)
93 plt.ylim([-0.05, 1.05])
94 plt.grid(True)
95 plt.tight_layout()
```

Example 3.5

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4
5  Example 3.5 k-fold cross-validationStochastic for Gradient Descent (SDG)
6  using the MINST data set
7
8  Developed for Machine Learning for Mechanical Engineers at the University of
9  South Carolina
10
11  @author: austin_downey
12  """
13
14  import IPython as IP
15  IP.get_ipython().magic('reset -sf')
16
17  import numpy as np
18  import scipy as sp
19  import pandas as pd
20  from scipy import fftpack, signal # have to add
21  import matplotlib as mpl
22  import matplotlib.pyplot as plt
23  import sklearn as sk
24  from sklearn import linear_model
25  from sklearn import pipeline
26  from sklearn import datasets
27
28  cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
29  plt.close('all')
30
31
32  """ Load your data
33
34  # Fetch the MNIST dataset from openml
35  mnist = sk.datasets.fetch_openml('mnist_784',as_frame=False)
36  X = mnist['data'] # load the data
37  Y = np.asarray(mnist['target'],dtype=int) # load the target
38
39  # Split the data set up into a training and testing data set
40  X_train = X[0:60000,:]
41  X_test = X[60000:,:]
42  Y_train = Y[0:60000]
43  Y_test = Y[60000:]
44
45  """ Train a Stochastic Gradient Descent classifier
46
47  # Extract a subset for our "5-detector".
48  Y_train_5 = (Y_train == 5)
49  Y_test_5 = (Y_test == 5)
50
51  sgd_clf = sk.linear_model.SGDClassifier()
52
53  """ Build a 5-detector several times so see the variation in metrics returned by SGD
54
55  # Solve the model multiple times to see the variations
56  for i in range(5):
57
58      # Train the model and return the predictions made by SGD
59      sgd_clf.fit(X_train, Y_train_5)
60      Y_train_pred = sgd_clf.predict(X_train)
61
62      # Use SK learn model to return metrics
63      accuracy = sk.metrics.accuracy_score(Y_train_5, Y_train_pred)
64      precision = sk.metrics.precision_score(Y_train_5, Y_train_pred)
65      recall = sk.metrics.recall_score(Y_train_5, Y_train_pred)
66      print('accuracy is '+str(np.round(accuracy,4))+'; precision is '+str(np.round(
precision,4))+

```

```
67         '; recall is '+str(np.round(recall,4))
68
69
70     """ Build a 5-detector using k-fold cross-validation
71
72     # From our example we can see quite a variation in results for a model, making it
73     # hard to select the proper model. However, k-fold cross-validation can help with this.
74
75
76     # make a prediction using the k-fold method to split up the data set. Again,
77     # solve the model multiple times to see the variations
78     for i in range(5):
79
80         # Train the model and return the predictions made by SGD
81         Y_train_pred = sk.model_selection.cross_val_predict(sgd_clf, X_train, Y_train_5, cv=3
82         )
83
84         # Use SK learn model to return metrics
85         accuracy = sk.metrics.accuracy_score(Y_train_5, Y_train_pred)
86         precision = sk.metrics.precision_score(Y_train_5, Y_train_pred)
87         recall = sk.metrics.recall_score(Y_train_5, Y_train_pred)
88         print('accuracy is '+str(np.round(accuracy,4))+'; precision is '+str(np.round(
89         precision,4))+
90         '; recall is '+str(np.round(recall,4)))
```

Example 3.6

```

1  """
2  Example 3.6 Multiclass Stochastic Gradient Descent (SDG) for the MNIST data set
3  @author: austin_downey
4  """
5
6  import IPython as IP
7  IP.get_ipython().run_line_magic('reset', '-sf')
8
9  import numpy as np
10 import matplotlib.pyplot as plt
11 import sklearn as sk
12 import time as time
13 from sklearn import linear_model
14 from sklearn import pipeline
15 from sklearn import datasets
16 from sklearn import multiclass
17
18 cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
19 plt.close('all')
20
21 """ Load your data
22
23 # Fetch the MNIST dataset from openml
24 mnist = sk.datasets.fetch_openml('mnist_784',as_frame=False,parser='auto')
25 X = np.asarray(mnist['data']) # load the data
26 Y = np.asarray(mnist['target'],dtype=int) # load the target
27
28 # Split the data set up into a training and testing data set
29 X_train = X[0:60000,:]
30 X_test = X[60000:,:]
31 Y_train = Y[0:60000]
32 Y_test = Y[60000:]
33
34 """ Train a Multiclass Stochastic Gradient Descent classifiers
35
36 # SK learn has a Multiclass and multilabel module as sk.multiclass. You can use
37 # this module to do one-vs-the-rest or one-vs-one classification.
38
39 # here we test a one-vs-rest classifier that uses Stochastic Gradient Descent
40 tt_1 = time.time()
41 ovr_clf = sk.multiclass.OneVsRestClassifier(sk.linear_model.SGDClassifier())
42 ovr_clf.fit(X_train, Y_train)
43 print('One-vs-Rest took '+str(time.time()-tt_1)+' seconds to train and execute')
44
45 # here we test a one-vs-one classifier that uses Stochastic Gradient Descent
46 tt_1 = time.time()
47 ovo_clf = sk.multiclass.OneVsOneClassifier(sk.linear_model.SGDClassifier())
48 ovo_clf.fit(X_train, Y_train)
49 print('One-vs-one took '+str(time.time()-tt_1)+' seconds to train and execute')
50
51 # Moreover, Scikit-Learn detects when you try to use a binary classification algorithm
52 # for
53 # a multiclass classification task, and it automatically runs OvA (except for SVM
54 # classifiers for which it uses OvO).
55 tt_1 = time.time()
56 multi_sgd_clf = sk.linear_model.SGDClassifier()
57 multi_sgd_clf.fit(X_train, Y_train) # y_train, not y_train_5
58 print('SK learns automated selection (OvA) took '+str(time.time()-tt_1)+' seconds to
59 train and execute')

```

Example 3.7

```

1  """
2  Example 3.7 Multiclass confusion matrix for the MNIST data set
3  @author: austin_downey
4  """
5
6  import IPython as IP
7  IP.get_ipython().run_line_magic('reset', '-sf')
8
9  import numpy as np
10 import matplotlib.pyplot as plt
11 import sklearn as sk
12
13 cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
14 plt.close('all')
15
16 """ Load your data
17
18 # Fetch the MNIST dataset from openml
19 mnist = sk.datasets.fetch_openml('mnist_784',as_frame=False,parser='auto')
20 X = np.asarray(mnist['data']) # load the data and convert to np array
21 Y = np.asarray(mnist['target'],dtype=int) # load the target
22
23 # Split the data set up into a training and testing data set
24 X_train = X[0:60000,:]
25 X_test = X[60000:,:]
26 Y_train = Y[0:60000]
27 Y_test = Y[60000:]
28
29 """ Confusion Matrix for a Multiclass classifier.
30
31 # Use the one-vs-one classifier that uses Stochastic Gradient Descent as this is
32 # faster for this specific data set
33 ovo_clf = sk.multiclass.OneVsOneClassifier(sk.linear_model.SGDClassifier())
34
35 # make a prediction for every case using the k-fold method.
36 Y_train_pred = sk.model_selection.cross_val_predict(ovo_clf, X_train, Y_train, cv=3)
37 conf_mx = sk.metrics.confusion_matrix(Y_train, Y_train_pred)
38 print(conf_mx)
39
40 # plot the results
41 fig = plt.figure(figsize=(4,4))
42 pos = plt.imshow(conf_mx) #, cmap=plt.cm.gray)
43 cbar = plt.colorbar(pos)
44 cbar.set_label('number of classified digits')
45 plt.ylabel('actual digit')
46 plt.xlabel('estimated digit')
47 plt.savefig('confusion_matrix',dpi=300)
48
49 # Normalize the confusion matrix by class size to compare error rates, not raw counts.
50 row_sums = conf_mx.sum(axis=1, keepdims=True)
51 norm_conf_mx = conf_mx / row_sums
52
53 # Next, we remove the high values along the diagonal. This is done by converting the
54 # confusion matrix to a float data type, and replacing everything on the diagonal with
55 # NaNs.
56 conf_mx_noise = np.asarray(norm_conf_mx,dtype=np.float32)
57 np.fill_diagonal(conf_mx_noise, np.NaN)
58
59 # plot the results but only consider the noise
60 fig = plt.figure(figsize=(4,4))
61 pos = plt.imshow(conf_mx_noise) #, cmap=plt.cm.gray)
62 cbar = plt.colorbar(pos)
63 cbar.set_label('normalized classification error')
64 plt.ylabel('actual digit')
65 plt.xlabel('estimated digit')
66 plt.savefig('confusion_matrix_error',dpi=300)

```

4 Regression-Based Classification

Certain algorithms that were first developed for regression can be adapted for classification, and some classifiers can be modified to predict continuous values. Converting a linear regressor into a classifier by adding a logistic link, for instance, retains the original coefficient vector, which makes it easy to see how each input feature influences the decision. As these dual-purpose models expose many of the hyper-parameters found in their regression versions, they often provide more opportunities for fine-tuning and clearer interpretability than algorithms designed purely for classification.

4.1 Logistic Regression

Logistic Regression, sometimes called logit regression, models the probability that a sample belongs to a specific class, such as the chance that an email is spam. When this probability exceeds 50%, the model predicts the positive class (label “1”); otherwise, it predicts the negative class (label “0”). It therefore acts as a binary classifier.

The predicted probability is

$$\hat{p} = h_{\theta}(X) = \sigma(\theta^T \cdot X). \quad (4.1)$$

Here \hat{p} is the estimated probability, and $\sigma(\cdot)$ is the sigmoid function, an S-shaped curve that maps any real number to the interval $(0, 1)$. The logistic function is defined in Equation 4.2 and illustrated in Figure 4.1. As before, $h_{\theta}(X)$ denotes the hypothesis that the input matrix X , augmented with a bias term, belongs to the positive class under the parameters θ .

$$\sigma(x) = \frac{1}{1 + e^{-x}}. \quad (4.2)$$

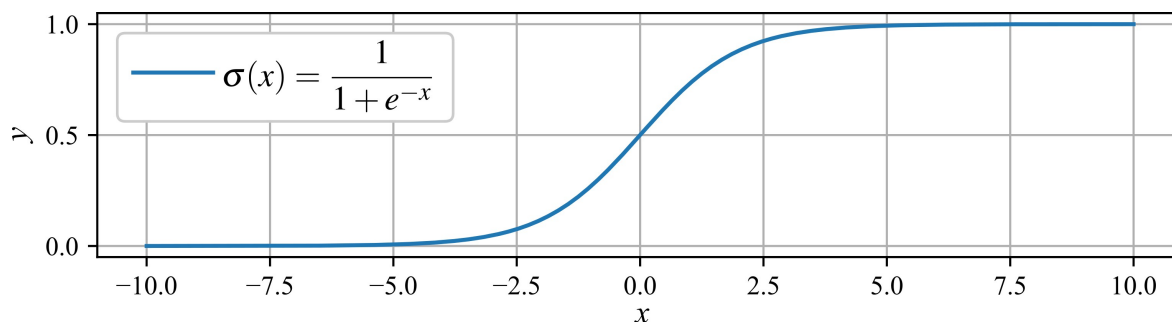


Figure 4.1: Sigmoid function that maps any real-valued input x to a value between 0 and 1.

Once the probability $\hat{p} = h_{\theta}(X)$ that an instance X belongs to the positive class has been estimated using Logistic Regression, the prediction (\hat{y}) can be made. \hat{y} is calculated as

$$\hat{y} = \begin{cases} 0 & \text{if } \hat{p} < 0.5, \\ 1 & \text{if } \hat{p} \geq 0.5. \end{cases} \quad (4.3)$$

Note that $\sigma(x) < 0.5$ when $x < 0$, and $\sigma(x) \geq 0.5$ when $x \geq 0$. Therefore, the Logistic Regression model predicts 0 if $\theta^T \cdot X$ is negative, and 1 if it is positive.

Now that we understand how logistic regression assigns probabilities and produces predictions, let's walk through a concise example that shows its training procedure and the associated cost function. Training seeks parameter values θ that give high predicted probabilities to positive examples ($y = 1$) and low probabilities to negative ones ($y = 0$). This aim is captured by the cost defined in equation 4.4, which evaluates a single training sample x . To achieve this, we require a cost function, such as

$$C(\theta) = \begin{cases} -\log(\hat{p}) & \text{if } y = 1, \\ -\log(1 - \hat{p}) & \text{if } y = 0. \end{cases} \quad (4.4)$$

which is plotted in 4.2.

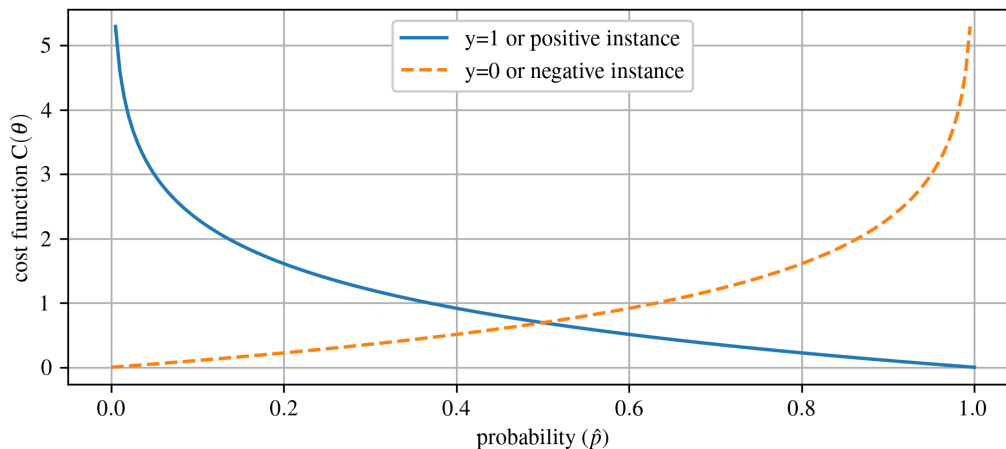


Figure 4.2: Cost function behavior for classification that heavily penalizes incorrect predictions.

The cost function in equation 4.4 behaves intuitively: $-\log(\hat{p})$ increases sharply as $\hat{p} \rightarrow 0$, so the loss is large when the model assigns a probability near 0 to a positive example. It likewise produces a high loss when the model predicts a probability close to 1 for a negative example. Conversely, because $-\log(\hat{p}) \rightarrow 0$ as $\hat{p} \rightarrow 1$, the loss becomes negligible when the predicted probability is near 1 for a positive instance or near 0 for a negative one, aligning with our expectations.

The overall cost, denoted $J(\theta)$, is the mean loss across all m training examples. This metric is commonly called the log-loss and written in equation 4.5:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(\hat{p}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{p}^{(i)}) \right]. \quad (4.5)$$

Because there is no closed-form solution for the parameters θ (embedded within \hat{p}), the minimum must be found with an iterative optimizer. The cost surface is convex, so gradient descent or another suitable algorithm will reach the global minimum provided the learning rate is reasonable and enough iterations are allowed. The gradient of the cost with respect to the j^{th} parameter θ_j is given in equation 4.6:

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m \left[\sigma(\theta^T X^{(i)}) - y^{(i)} \right] x_j^{(i)}. \quad (4.6)$$

This equation resembles the partial derivative used in gradient descent. For every training example, it finds the prediction error, multiplies it by the j^{th} feature value, and then averages these products over all m samples. With the resulting gradient vector of partial derivatives, you can update the parameters using the Batch Gradient Descent algorithm, completing the training of a Logistic Regression model. Stochastic Gradient Descent performs the update after each single example, while Mini-batch Gradient Descent updates the parameters after processing each mini-batch.

Review 4.1 Iris Flower Dataset

The Iris dataset, collected by botanist Edgar Anderson in 1935, lists sepal length, sepal width, petal length, and petal width for 150 flowers drawn equally from three species: Iris-Setosa, Iris-Versicolor, and Iris-Virginica (refer to Figure 4.3)a. Statistician Ronald Fisher analysed these measurements in a 1936 study on linear discriminant analysis, and the data have since become a standard benchmark in statistics and machine learning because they are small, clean, and perfectly balanced, making them ideal for testing classification algorithms and visualisation techniques.

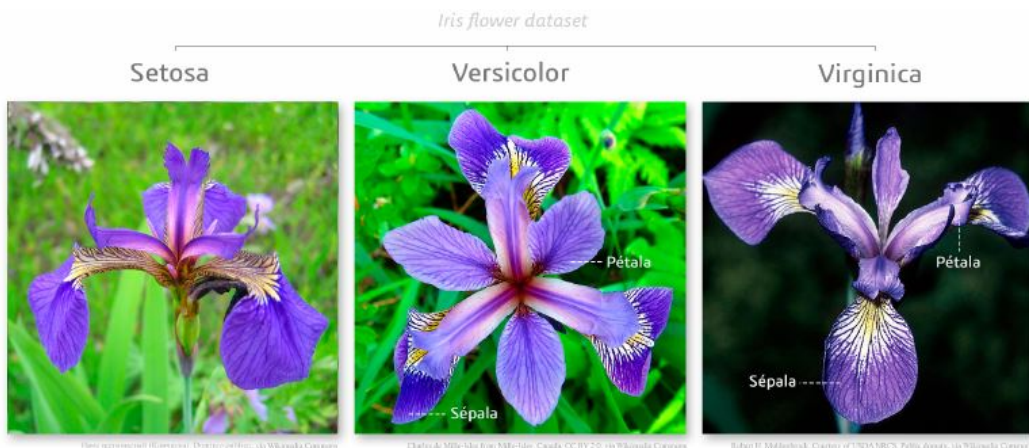


Figure 4.3: Flowers representing three species of iris plants ^a

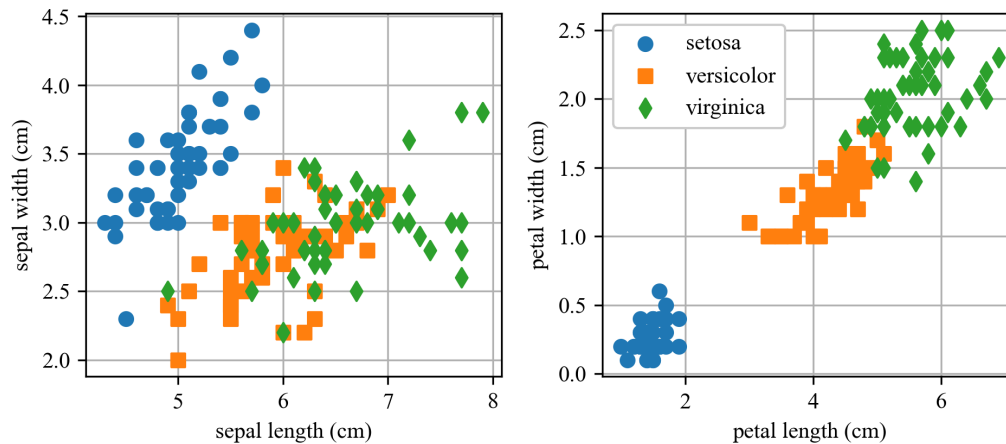


Figure 4.4: Iris dataset scatterplot showing sepal length vs. sepal width (left) and petal length vs. petal width (right).

^aDiego Mariano, CC BY-SA 4.0 <<https://creativecommons.org/licenses/by-sa/4.0/>>, via Wikimedia Commons

Example 4.1 Iris Dataset Exploration

This example introduces the Iris dataset (Figure 4.3), originally compiled by Ronald Fisher. It demonstrates how to load the dataset, access feature and label information, and visualize the relationships between sepal and petal measurements across the three iris species.

4.1.1 1-D Decision Boundaries

The petal width of Iris-virginica flowers usually lies between 1.4 cm and 2.5 cm, whereas the other two species range from 0.1 cm to 1.8 cm. Because these intervals overlap, the classifier's certainty changes across the scale. Above roughly 2 cm, the model is confident that any given sample is Iris-virginica. Moreover, below 1 cm it is equally sure the flower is not Iris-virginica (This means the model predicts high probability for the class "Not Iris-Virginica"). In between the limits of the classes, the model is unsure. When you call `predict()` instead of `predict_proba()`, it outputs the most likely class, creating a decision boundary around 1.6 cm where both class probabilities reach 50 percent. Thus, flowers with petal widths greater than 1.6 cm lead to a prediction of Iris-virginica, while flowers with smaller petal widths yield the opposite prediction even if confidence is low.

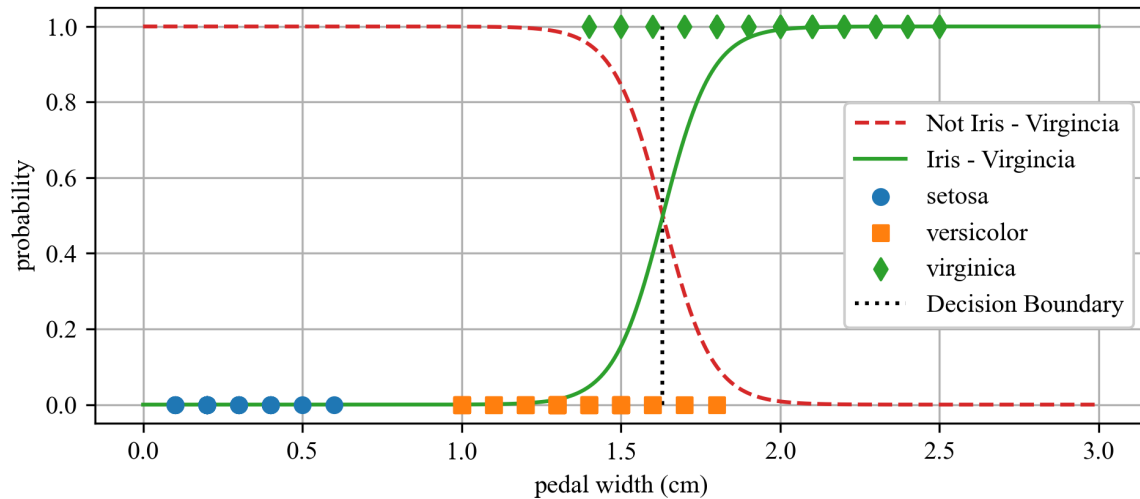


Figure 4.5: Decision boundary for the flowers of three Iris plant species with C set to $C = 10^{10}$.

Example 4.2 1D Logistic Regression

This example builds a logistic regression model to classify Iris-Virginica based on a single feature: petal width. It visualizes class probabilities and shows the decision boundary at the 50% threshold.

4.1.2 2-D Decision Boundaries

Figure 4.6 plots petal width against petal length for the Iris samples. After training, the logistic regression model assigns each point the probability that the flower is *Irisvirginica*. The dashed line marks where this probability equals 50%, forming the decision boundary, which is linear. The parallel contour lines show equal-probability levels from 10% to 90%. Points beyond the top-right line have a greater than 90% probability of being classified as Iris-Virginica by the model.

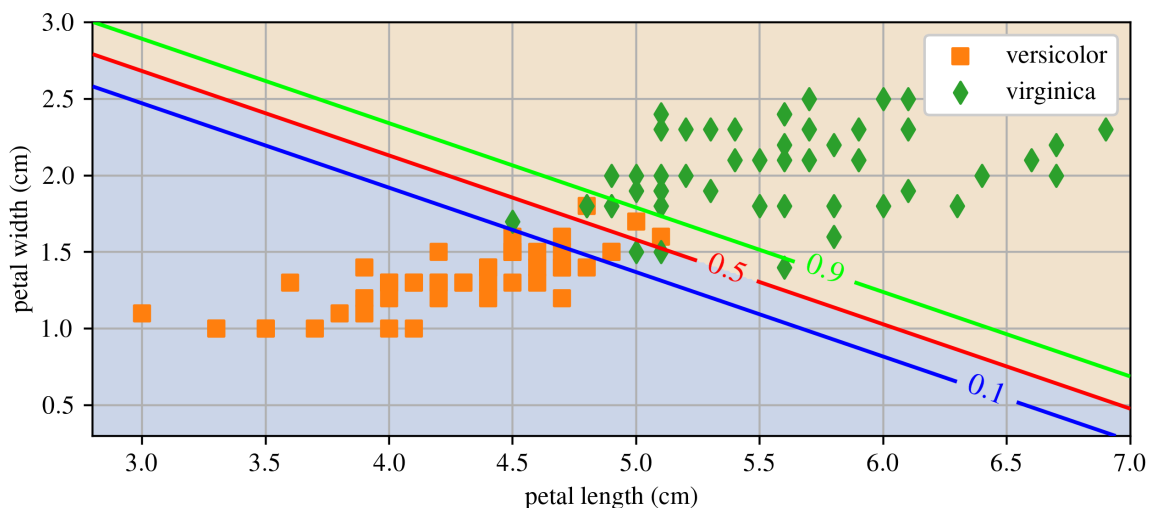


Figure 4.6: A 2D decision boundary for the Iris dataset.

Example 4.3 2D Decision Boundary

This example uses logistic regression to classify the Iris-Virginica species based on petal length and width. A 2D decision boundary is visualized in the petal feature space, with prediction regions and probability contours illustrating classifier confidence.

4.2 Softmax Regression

Logistic Regression can be generalised to handle many classes in a single model, so there is no need to train and merge several binary classifiers. This multiclass version is called Softmax Regression, or Multinomial Logistic Regression. Softmax Regression provides a straightforward way to perform regression-based classification when more than two classes are present. For an input vector \mathbf{x} , the model first computes a score $s_k(\mathbf{x})$ for every class k , then turns these scores into class probabilities with the softmax (normalised exponential) function. The score is obtained exactly as in linear regression:

$$s_k(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\theta}^{(k)}. \quad (4.7)$$

Here each class has its own parameter vector $\boldsymbol{\theta}^{(k)}$, and the collection of these vectors is usually stored as the rows of a parameter matrix Θ .

After the model has calculated a score for each class given an input \mathbf{x} , it converts these scores to probabilities with the softmax function. For class k the predicted probability is

$$\hat{p}_k = \sigma(\mathbf{s}(\mathbf{x}))_k = \frac{\exp(s_k(\mathbf{x}))}{\sum_{j=1}^K \exp(s_j(\mathbf{x}))}. \quad (4.8)$$

In this expression, $\mathbf{s}(\mathbf{x})$ is the vector of class scores for the input, $\sigma(\mathbf{s}(\mathbf{x}))_k$ is the probability that \mathbf{x} belongs to class k , and K is the total number of classes. In short, equation 4.8 exponentiates each score and then normalises the results by dividing by the sum of all exponentials so the probabilities sum to one.

Like logistic regression, a Softmax Regression model assigns an input to the class whose predicted probability is largest, which coincides with the class that has the highest score:

$$\hat{y} = \underset{k}{\operatorname{argmax}} \sigma(\mathbf{s}(\mathbf{x}))_k = \underset{k}{\operatorname{argmax}} s_k(\mathbf{x}) = \underset{k}{\operatorname{argmax}} ((\boldsymbol{\theta}^{(k)})^T \cdot \mathbf{x}). \quad (4.9)$$

NOTE

The argmax operator returns the argument that maximises a function. In this setting, it yields the index k for which the estimated probability $\sigma(\mathbf{s}(\mathbf{x}))_k$ attains its largest value.

With probability estimation and prediction established, we next consider training. The task is to learn parameters that place a large probability on the correct class and correspondingly small probabilities on all others. This objective is met by minimising the cost in equation 4.10, known as the cross entropy, which heavily penalises the model when it assigns a low probability to the true

label. Cross entropy is widely used to measure how well predicted class probabilities agree with the actual classes. The cost function is

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(\hat{p}_k^{(i)}). \quad (4.10)$$

NOTE

In this expression, $y_k^{(i)} = 1$ when the i^{th} sample's true label is class k , and $y_k^{(i)} = 0$ otherwise.

When considering only two classes ($K = 2$), it's important to highlight that this cost function aligns with the Logistic Regression's cost function, commonly referred to as log loss (refer to Equation 4.5).

The gradient of the cross-entropy cost with respect to $\theta^{(k)}$ is

$$\nabla_{\theta^{(k)}} J(\Theta) = \frac{1}{m} \sum_{i=1}^m (\hat{p}_k^{(i)} - y_k^{(i)}) X^{(i)}. \quad (4.11)$$

By computing this vector for each class, you obtain the full gradient, which can then be fed to Gradient Descent or another optimizer to find the parameter matrix Θ that minimises the cost.

Applying Softmax Regression to the three-class iris problem in Scikit-Learn is straightforward. The `LogisticRegression` estimator normally uses one-versus-all when more than two classes are present, but setting `multi_class='multinomial'` activates true Softmax learning. Choose a solver that supports this option, such as `lbfgs` (see the library documentation). The model applies ℓ_2 regularisation by default, governed by the hyperparameter C . After training, a flower with 5 cm long and 2 cm wide petals is classified as Iris-Virginica with probability 94.2%, while the probability of Iris versicolor is 5.8%.

NOTE

The Softmax Regression classifier is multiclass, not multioutput. As such, Softmax Regression can only predict one class at a time, so it works for problems where each input belongs to exactly one category-like classifying an email as spam, promotions, or updates. It can't be used for cases where multiple labels may apply, such as tagging a news article with topics like politics, economics, and technology all at once.

Figure 4.7 displays the decision regions, shaded with background colours to mark each class. The borders separating any two classes are straight lines. The plot also includes contour curves for the predicted probability of the Iris-Versicolor class. At the point where all three borders meet, every class receives a probability of 33%, so the selected class can have a confidence below 50%.

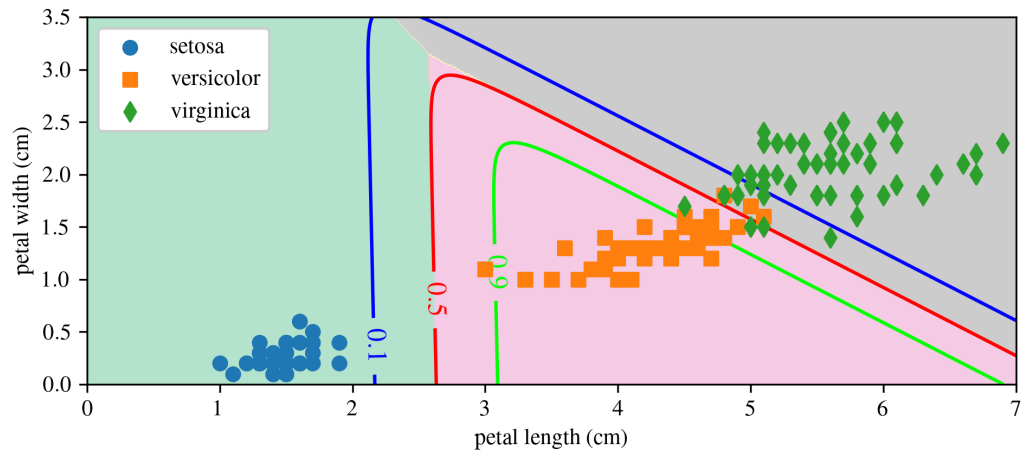


Figure 4.7: Softmax classification for the three iris plant species.

Example 4.4 Softmax Classification

This example applies softmax regression to classify all three iris species using petal length and width. It builds a multinomial logistic regression model and visualizes the decision boundaries along with confidence contours over the 2D petal feature space.

4.3 Examples

Example 4.1

```
1  """
2  Example 4.1 Introduction to the IRIS data set
3  @author: austin_downey
4  """
5
6  import IPython as IP
7  IP.get_ipython().run_line_magic('reset', '-sf')
8
9  import matplotlib.pyplot as plt
10 import sklearn as sk
11
12 cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
13 plt.close('all')
14
15
16 %% Load your data
17
18 # We will use the Iris data set. This dataset was created by biologist Ronald
19 # Fisher in his 1936 paper "The use of multiple measurements in taxonomic
20 # problems" as an example of linear discriminant analysis
21
22 iris = sk.datasets.load_iris()
23
24 # for simplicity, extract some of the data sets
25 X = iris['data'] # this contains the length of the pedals and sepals
26 Y = iris['target'] # contains what type of flower it is
27 Y_names = iris['target_names'] # contains the name that aligns with the type of the
   flower
28 feature_names = iris['feature_names'] # the names of the features
29
30 # plot the Sepal data
31 plt.figure(figsize=(6.5,3))
32 plt.subplot(121)
33 plt.grid(True)
34 plt.scatter(X[Y==0,0],X[Y==0,1],marker='o')
35 plt.scatter(X[Y==1,0],X[Y==1,1],marker='s')
36 plt.scatter(X[Y==2,0],X[Y==2,1],marker='d')
37 plt.xlabel(feature_names[0])
38 plt.ylabel(feature_names[1])
39
40
41 plt.subplot(122)
42 plt.grid(True)
43 plt.scatter(X[Y==0,2],X[Y==0,3],marker='o',label=Y_names[0])
44 plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1])
45 plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2])
46 plt.xlabel(feature_names[2])
47 plt.ylabel(feature_names[3])
48 plt.legend(framealpha=1)
49 plt.tight_layout()
50
```

Example 4.2

```

1  """
2  Example 4.2 1D Decision boundary for the Iris dataset
3  @author: austin_downey
4  """
5
6  import IPython as IP
7  IP.get_ipython().run_line_magic('reset', '-sf')
8
9  import numpy as np
10 import matplotlib.pyplot as plt
11 import sklearn as sk
12
13
14 cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
15 plt.close('all')
16
17
18 """ Load your data
19
20 # We will use the Iris data set. This dataset was created by biologist Ronald
21 # Fisher in his 1936 paper "The use of multiple measurements in taxonomic
22 # problems" as an example of linear discriminant analysis
23
24 iris = sk.datasets.load_iris()
25
26 # for simplicity, extract some of the data sets
27 X = iris['data'] # this contains the length of the pedals and sepals
28 Y = iris['target'] # contains what type of flower it is
29 Y_names = iris['target_names'] # contains the name that aligns with the type of the
30 flower
31 feature_names = iris['feature_names'] # the names of the features
32
33 # plot the Sepal data
34 plt.figure(figsize=(6.5,3))
35 plt.subplot(121)
36 plt.grid(True)
37 plt.scatter(X[Y==0,0],X[Y==0,1],marker='o')
38 plt.scatter(X[Y==1,0],X[Y==1,1],marker='s')
39 plt.scatter(X[Y==2,0],X[Y==2,1],marker='d')
40 plt.xlabel(feature_names[0])
41 plt.ylabel(feature_names[1])
42
43 plt.subplot(122)
44 plt.grid(True)
45 plt.scatter(X[Y==0,2],X[Y==0,3],marker='o',label=Y_names[0])
46 plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1])
47 plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2])
48 plt.xlabel(feature_names[2])
49 plt.ylabel(feature_names[3])
50 plt.legend(framealpha=1)
51 plt.tight_layout()
52
53
54 """ Train a Logistic Regression model
55
56 # define the features (X) and the output (Y)
57 X_pedal = iris["data"][:, 3:] # consider just the petal width
58 y_pedal = iris["target"] == 2 # 1 if Iris-Virginica, else 0
59
60 # Build the logistic Regression model and train it.
61 log_reg = sk.linear_model.LogisticRegression(C=1)
62 log_reg.fit(X_pedal, y_pedal)
63 # Note: The hyper-parameter controlling the regularization strength of a
64 # Scikit-Learn LogisticRegression model is not alpha (as in other linear models),
65 # but its inverse: C. The higher the value of C, the less the model is regularized.
66

```

```
67 # Build a range of the feature (X) to predict over. Here we just consider pedal width.
68 X_new = np.linspace(0, 3, 1000)
69 X_new = np.expand_dims(X_new, axis=1)
70
71 # Use the Logistic Regression Model to predict the pedal type based on pedal width
72 y_proba = log_reg.predict_proba(X_new)
73
74 # plot the probability plots
75 plt.figure(figsize=(6.5,3))
76 plt.grid(True)
77
78 # plot the data used for training, set at 0 and 1
79 plt.scatter(X[Y==0,3],np.zeros(50),marker='o',label=Y_names[0],zorder=10)
80 plt.scatter(X[Y==1,3],np.zeros(50),marker='s',label=Y_names[1],zorder=10)
81 plt.scatter(X[Y==2,3],np.ones(50),marker='d',label=Y_names[2],zorder=10)
82
83 # plot the probability
84 plt.plot(X_new,y_proba[:,0], '--',color=cc[3],label='Not Iris - Virgincia')
85 plt.plot(X_new,y_proba[:,1],color=cc[2], label='Iris - Virgincia')
86
87 # find and plot the 50% decision boundary
88 x_at_50 = X_new[np.argmin(np.abs(y_proba[:,0] - 0.5)))]
89 plt.vlines(x_at_50,0,1,color='k',linestyles=':',label='Decision Boundary')
90
91 plt.xlabel('pedal width (cm)')
92 plt.ylabel('probability')
93 plt.legend(framealpha=1)
94 plt.tight_layout()
95
96 # make predictions on the model trained on the data.
97 log_reg.predict([[1.7]])
```

Example 4.3

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Example 4.3 2D Decision boundary for the Iris dataset
5
6  Developed for Machine Learning for Mechanical Engineers at the University of
7  South Carolina
8
9  @author: austin_downey
10 """
11
12 import IPython as IP
13 IP.get_ipython().run_line_magic('reset', '-sf')
14
15
16 import numpy as np
17 import matplotlib.pyplot as plt
18 import sklearn as sk
19
20
21 cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
22 plt.close('all')
23
24
25 """ Load your data
26
27 # We will use the Iris data set. This dataset was created by biologist Ronald
28 # Fisher in his 1936 paper "The use of multiple measurements in taxonomic
29 # problems" as an example of linear discriminant analysis
30
31 iris = sk.datasets.load_iris()
32
33 # for simplicity, extract some of the data sets
34 X = iris['data'] # this contains the length of the petals and sepals
35 Y = iris['target'] # contains what type of flower it is
36 Y_names = iris['target_names'] # contains the name that aligns with the type of the
37 flower
38 feature_names = iris['feature_names'] # the names of the features
39
40 # plot the Sepal data
41 plt.figure(figsize=(6.5,3))
42 plt.subplot(121)
43 plt.grid(True)
44 plt.scatter(X[Y==0,0],X[Y==0,1],marker='o',zorder=10)
45 plt.scatter(X[Y==1,0],X[Y==1,1],marker='s',zorder=10)
46 plt.scatter(X[Y==2,0],X[Y==2,1],marker='d',zorder=10)
47 plt.xlabel(feature_names[0])
48 plt.ylabel(feature_names[1])
49
50 plt.subplot(122)
51 plt.grid(True)
52 plt.scatter(X[Y==0,2],X[Y==0,3],marker='o',label=Y_names[0],zorder=10)
53 plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1],zorder=10)
54 plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2],zorder=10)
55 plt.xlabel(feature_names[2])
56 plt.ylabel(feature_names[3])
57 plt.legend(framealpha=1)
58 plt.tight_layout()
59
60
61 """ plot the Linear decision boundary in 2D "Petal" space
62
63 # build the training and target set.
64 X_train = X[:, (2, 3)] # petal length, petal width
65 y_train = Y == 2
66

```

```
67 # build the Logistic Regression model
68 log_reg = sk.linear_model.LogisticRegression(C=10**10)
69 # Note: The hyper-parameter controlling the regularization strength of a Scikit-Learn
70 # LogisticRegression model is not alpha (as in other linear models), but its
71 # inverse: C. The higher the value of C, the less the model is regularized.
72
73 # train the Logistic Regression model
74 log_reg.fit(X_train, y_train)
75
76 # build the x values for the predictions over the entire "petal space"
77 x_grid, y_grid = np.meshgrid(
78     np.linspace(2.8, 7, 500),
79     np.linspace(0.3, 3, 200),
80 )
81 X_new = np.vstack((x_grid.reshape(-1), y_grid.reshape(-1))).T # build a vector format of
the mesh grid
82
83 # predict on the vectorized format
84 y_predict = log_reg.predict(X_new)
85 y_proba = log_reg.predict_proba(X_new)
86
87 # convert back to meshgrid shape for plotting
88 zz_predict = y_predict.reshape(x_grid.shape)
89 zz_proba = y_proba[:, 1].reshape(x_grid.shape)
90
91 # plot the 2D "petal space"
92 plt.figure(figsize=(6.5,3))
93 plt.grid(True)
94 plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',color=cc[1],label=Y_names[1],zorder=10)
95 plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',color=cc[2],label=Y_names[2],zorder=10)
96 plt.contourf(x_grid, y_grid, zz_predict, cmap='Pastel2')
97 contour = plt.contour(x_grid, y_grid, zz_proba, [0.100,0.5,0.900],cmap=plt.cm.brg)
98 plt.clabel(contour, inline=1) # add the labels to the plot
99 plt.xlabel(feature_names[2])
100 plt.ylabel(feature_names[3])
101 plt.legend()
102 plt.tight_layout()
103
104
105
```

Example 4.4

```

1  """
2  Example 4.4 Softmax decision boundary for the Iris dataset
3  @author: austin_downey
4  """
5
6  import IPython as IP
7  IP.get_ipython().magic('reset -sf')
8
9  import numpy as np
10 import matplotlib.pyplot as plt
11 import sklearn as sk
12
13
14 cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
15 plt.close('all')
16
17
18 """ Load your data
19
20 # We will use the Iris data set. This dataset was created by biologist Ronald
21 # Fisher in his 1936 paper "The use of multiple measurements in taxonomic
22 # problems" as an example of linear discriminant analysis
23 iris = sk.datasets.load_iris()
24
25 # for simplicity, extract some of the data sets
26 X = iris['data'] # this contains the length of the petals and sepals
27 Y = iris['target'] # contains what type of flower it is
28 Y_names = iris['target_names'] # contains the name that aligns with the type of the
29 flower
30 feature_names = iris['feature_names'] # the names of the features
31
32 # plot the Sepal data
33 plt.figure(figsize=(6.5,3))
34 plt.subplot(121)
35 plt.grid(True)
36 plt.scatter(X[Y==0,0],X[Y==0,1],marker='o',zorder=10)
37 plt.scatter(X[Y==1,0],X[Y==1,1],marker='s',zorder=10)
38 plt.scatter(X[Y==2,0],X[Y==2,1],marker='d',zorder=10)
39 plt.xlabel(feature_names[0])
40 plt.ylabel(feature_names[1])
41
42 plt.subplot(122)
43 plt.grid(True)
44 plt.scatter(X[Y==0,2],X[Y==0,3],marker='o',label=Y_names[0],zorder=10)
45 plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1],zorder=10)
46 plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2],zorder=10)
47 plt.xlabel(feature_names[2])
48 plt.ylabel(feature_names[3])
49 plt.legend(framealpha=1)
50 plt.tight_layout()
51
52 """ Softmax Regression
53
54 # build the training and target set.
55 X_train = X[:, (2, 3)] # petal length, petal width
56 y_train = Y
57
58 # build and train the softmax model
59 softmax_reg = sk.linear_model.LogisticRegression(multi_class="multinomial",
60 solver="lbfgs", C=10)
61 softmax_reg.fit(X_train, y_train)
62
63 # build the x values for the predictions over the entire "petal space"
64 x_grid, y_grid = np.meshgrid(
65 np.linspace(0, 7, 500),
66 np.linspace(0, 4, 200),

```

```
67     )
68     X_new = np.vstack((x_grid.reshape(-1), y_grid.reshape(-1))).T # build a vector format of
        the mesh grid
69
70     # predict on the vectorized format
71     y_predict = softmax_reg.predict(X_new)
72     y_proba = softmax_reg.predict_proba(X_new)
73
74
75     # convert back to meshgrid shape for plotting
76     zz_predict = y_predict.reshape(x_grid.shape)
77     zz_proba = y_proba[:, 1].reshape(x_grid.shape) # the selected column selects the
        probability that the data falls within this class.
78
79     # plot the 2D "petal space"
80     plt.figure(figsize=(6.5, 4))
81     plt.scatter(X[Y==0,2],X[Y==0,3],marker='o',label=Y_names[0],zorder=10)
82     plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1],zorder=10)
83     plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2],zorder=10)
84     plt.contourf(x_grid, y_grid, zz_predict, cmap='Pastel2')
85     contour = plt.contour(x_grid, y_grid, zz_proba, [0.100,0.5,0.900], cmap=plt.cm.brg)
86     plt.clabel(contour, inline=1)
87     plt.xlabel(feature_names[2])
88     plt.ylabel(feature_names[3])
89     plt.legend()
90     plt.tight_layout()
```

5 Decision Trees

Decision trees are flexible learners that work well for classification, regression, and multi-output problems. A single tree can capture intricate relationships in the data with minimal preprocessing. Each tree also acts as the core building block of ensemble methods such as Random Forests, which are among the most reliable models in modern practice. One limitation is size: a fully grown tree can become quite large, so pruning or depth limits are often applied to control memory use and reduce overfitting.

In this chapter, we will explore the essentials of Decision Trees, starting with their training, visualization, and prediction processes. We will then look into the CART (Classification and Regression Trees) training algorithm, which Scikit-Learn utilizes for constructing Decision Trees. Additionally, we will examine how to regulate the complexity of Decision Trees and adapt them for regression tasks. The chapter concludes by addressing some inherent limitations of Decision Trees.

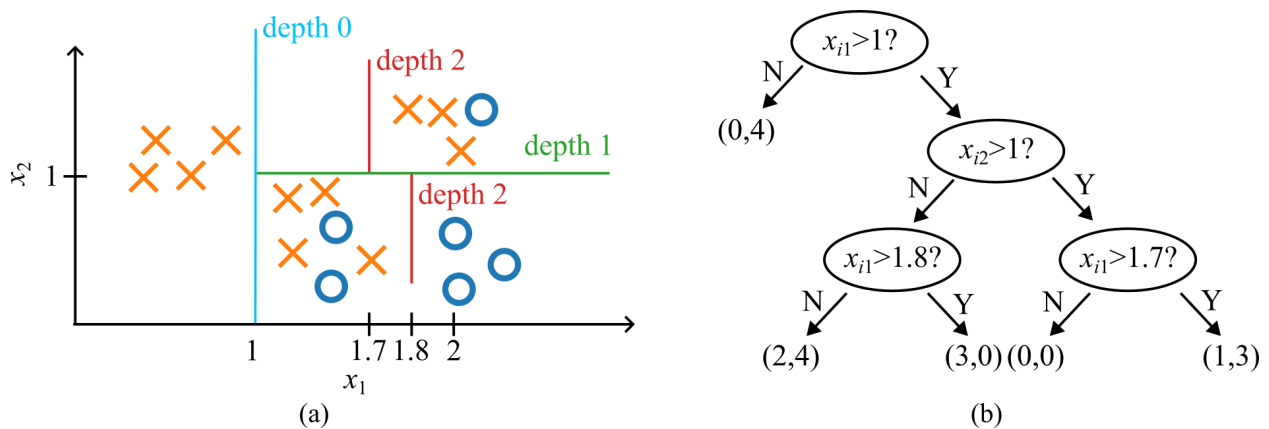


Figure 5.1: The basic connect of a decision tree, showing (a) how a decision tree is built, and (b) the developed decision tree.

5.1 Decision Tree Classification

Figure 5.2 shows the regions formed by the decision tree. The solid vertical line marks the root split at depth 0, where petal length equals 2.45 cm. All samples on the left of this line belong solely to Iris setosa, so no further division is needed there. The right half still mixes species, so the depth-1 right child splits again at petal width 1.75 cm, indicated by the dashed line. With `max_depth=2` the tree stops at this point. If `max_depth` were increased to 3, the two depth-2 leaves would introduce additional boundaries, drawn as dotted lines.

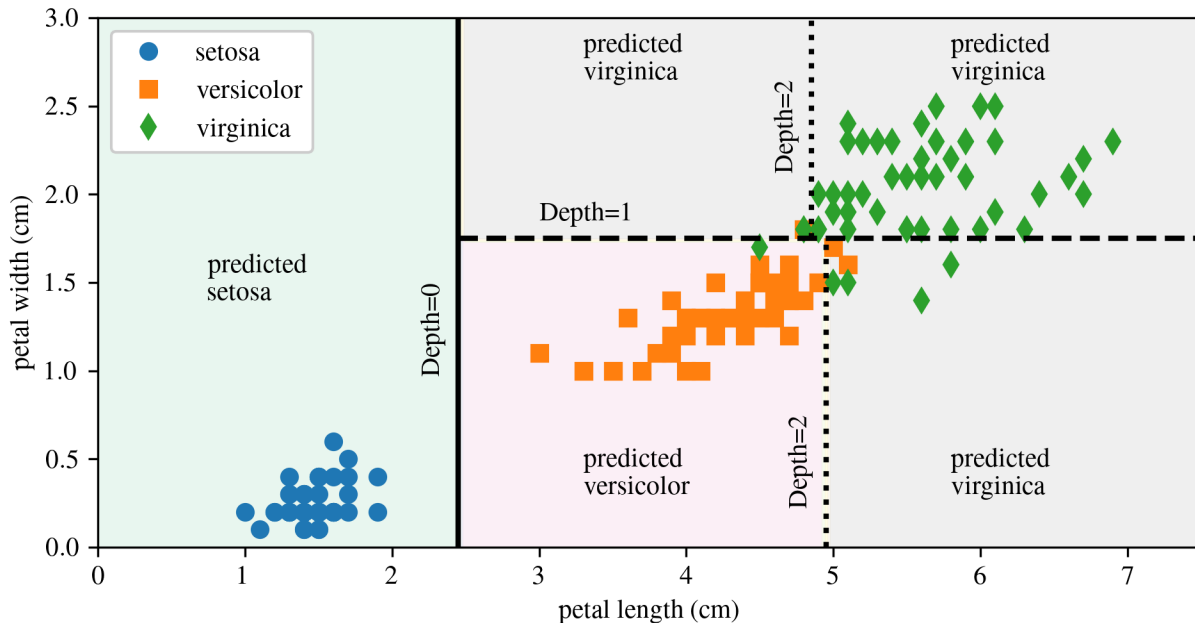


Figure 5.2: Decision Tree decision boundaries.

A decision tree for the Iris Dataset looks like Figure 5.3. Let us examine how the Decision Tree processes predictions. Suppose you come across an iris and need to classify it. Start at the root node (depth=0), which asks whether the flower's petal length is less than 2.45 cm. If the answer is yes, move to the root's left child (depth=1). This child is a leaf, so no further questions follow. The class stored in that leaf is Iris-setosa, and the tree therefore labels the sample as setosa.

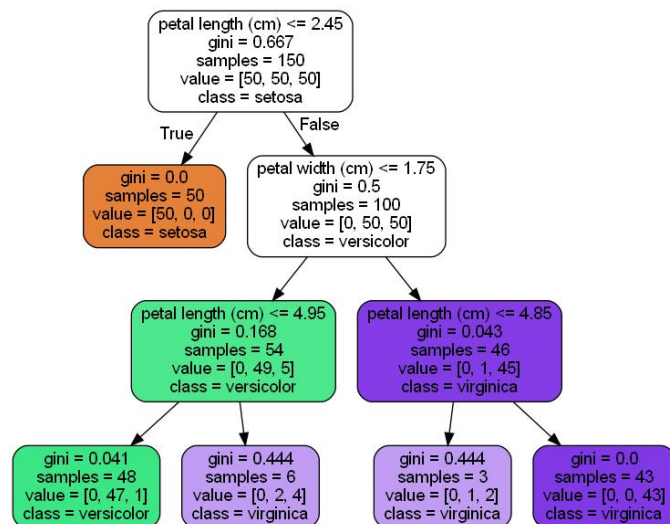


Figure 5.3: Decision tree for the Iris Dataset trained using the CART algorithms.

In Figure 5.3, each node reports the split criterion, Gini impurity, sample count, and class distribution. Fill colors denote the predicted species; orange for Setosa, green for Versicolor, and

purple for Virginica. Shade intensity conveys confidence: darker hues indicate purer (more certain) leaves, whereas white represents complete uncertainty. Consider a second iris whose petal length is greater than 2.45 cm. From the root you move to its right child (depth 1). This internal node asks a new question: is the petal width less than 1.75 cm? If yes, the flower is classified as Iris-versicolor (depth=2, left leaf). If no, it is labeled Iris-virginica (depth=2, right leaf). The decision path is clear and easy to follow.

The attributes of a node include:

- *class*: the class label stored in the leaf.
- *samples*: the number of training instances that reach the node. For example, 100 flowers have petal length >2.45 cm (depth1, right); among them, 54 also have petal width < 1.75 cm (depth2, left).
- *value*: a three-element vector giving the count of instances from each species at the node. The bottom-right leaf, for instance, contains 0 Iris-setosa, 1 Iris-versicolor, and 45 Iris-virginica.
- *gini*: the Gini impurity. A node is *pure* when $gini = 0$, meaning all samples belong to the same class. The depth-1 left node, which holds only Iris-setosa, is pure.

The Gini impurity G_i for the i^{th} node is computed using

$$G_i = 1 - \sum_{k=1}^n p_{i,k}^2, \quad (5.1)$$

where $p_{i,k}$ is the proportion of class k instances among the training instances at the i^{th} node. For instance, the Gini score for the depth-2 left node is calculated as follows: $1 - (0/54)^2 - (49/54)^2 - (5/54)^2 \approx 0.168$. We will discuss an alternative impurity measure later.

5.1.1 Class Probability Estimation

A Decision Tree can estimate the likelihood that a sample belongs to class k . The algorithm traces the sample's path to its leaf node and then reports the fraction of training instances of class k found in that leaf. Consider a flower with petals 5cm long and 1.5cm wide. This sample reaches the depth-2 left leaf, which contains 54 training flowers: 0 *Iris-setosa*, 49 *Iris-versicolor*, and 5 *Iris-virginica*. The tree therefore assigns probabilities of 0%, 90.7%, and 9.3% to the three classes, respectively. When asked for a class label, it selects *Iris-versicolor*, the class with the highest probability.

Interestingly, the same estimated probabilities apply throughout the bottom-right rectangle of Figure 5.2, even in a scenario where the petal dimensions are 6 cm by 1.5 cm where one might intuitively expect an Iris-Virginica classification.

Example 5.1 Decision Tree Classifier

This example uses Scikit-Learn’s DecisionTreeClassifier to classify iris species based on petal features. The decision process is visualized using Graphviz, showing the structure of learned rules and how decisions are made based on feature thresholds. This online viewer is the easiest way to do that <https://dreampuf.github.io/GraphvizOnline/>

5.2 The CART Training Algorithm

Scikit-Learn trains decision trees with the Classification and Regression Tree (CART) algorithm^a, a process often called “growing” the tree. At each node the algorithm splits the local training set into two parts by selecting a feature k and a threshold t_k , for example “petal length ≤ 2.45 cm”. It evaluates every candidate pair (k, t_k) and chooses the one that yields the purest child nodes, judging purity by a weighted combination of each child’s impurity measure and sample count.

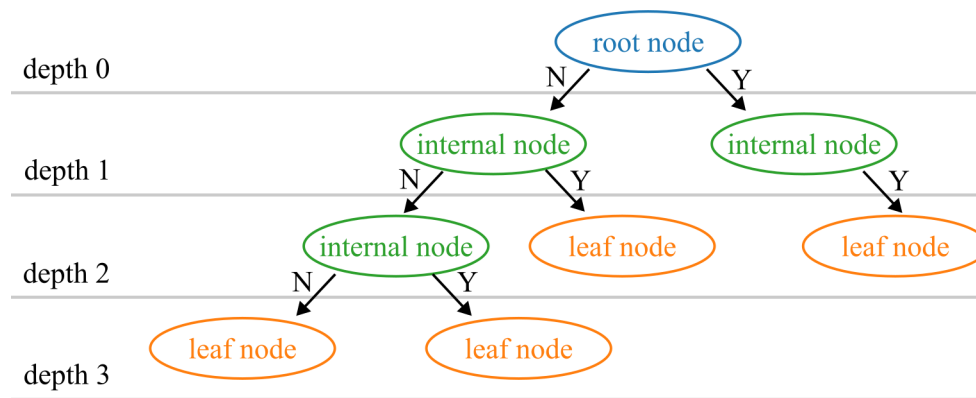


Figure 5.4: Illustration of the CART algorithm process.

The objective function it aims to minimize during this process is

$$J(k, t_k) = \frac{m_{\text{left}}}{m} G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}}. \quad (5.2)$$

In this equation, $G_{\text{left/right}}$ measures the impurity of the left/right subset, and $m_{\text{left/right}}$ is the number of instances in the left/right subset. After the initial split, the algorithm continues to divide the resulting subsets and their subsequent divisions recursively. This recursive process halts when it reaches the predefined maximum depth (`max_depth`) set through a hyperparameter, or if no further impurity-reducing splits can be found.

Unlike linear models that assume a specific data structure, Decision Trees impose few assumptions on their initial data. If not appropriately constrained, a Decision Tree can intricately conform to the training data, leading to overfitting. This model type is described as nonparametric, not due to a lack of parameters, which it can have in abundance, but because the parameters are not predetermined before training, allowing the model’s structure to freely mirror the data intricacies. Conversely, a parametric model, like a linear model, has a fixed number of parameters, limiting its flexibility but also minimizing overfitting risks while potentially increasing underfitting risks.

^aBreiman, Leo. Classification and regression trees. Routledge, 2017

To mitigate overfitting in Decision Trees, it is essential to control the model's freedom during training through regularization. The regularization hyperparameters vary by the algorithm, but typically, the tree's maximum depth can be restricted. In Scikit-Learn, this is managed by the `max_depth` hyperparameter, which is unlimited by default. Lowering `max_depth` helps regularize the model, thereby reducing overfitting likelihood.

Other parameters in Scikit-Learn's `DecisionTreeClassifier` also influence the tree's structure:

- `min_samples_split`: The minimum number of samples required to split a node.
- `min_samples_leaf`: The minimum number of samples a leaf node must have.
- `min_weight_fraction_leaf`: Similar to `min_samples_leaf`, but expressed as a fraction of total weighted instances.
- `max_leaf_nodes`: The maximum number of leaf nodes.
- `max_features`: The maximum number of features evaluated for splitting at each node.

Adjusting these parameters by increasing `min_*` values or decreasing `max_*` values will help regularize the model.

Figure 5.5 illustrates two Decision Trees trained on the moons dataset: one on the left with default hyperparameters (unrestricted) and another on the right with `min_samples_leaf=4`. The left model appears to be overfitting, whereas the right model, with its restrictions, likely offers better generalization.

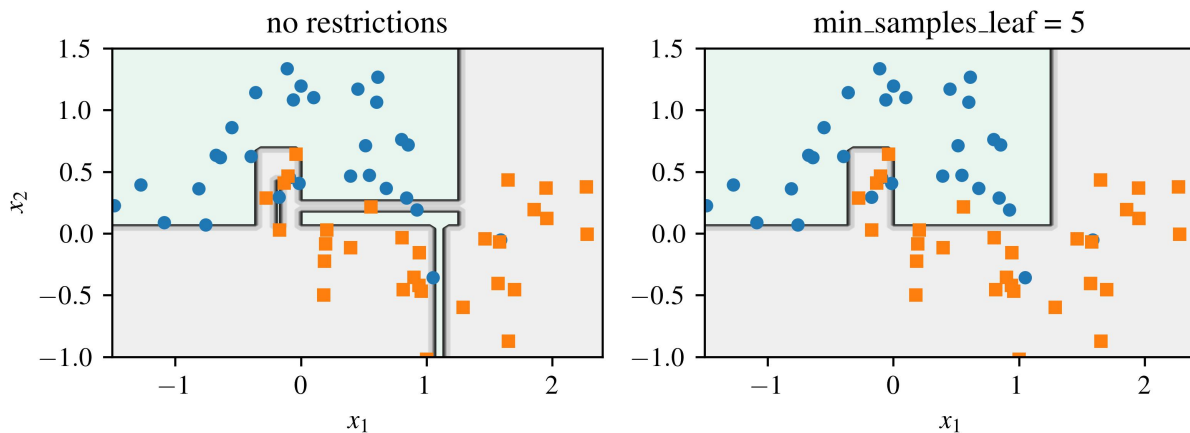


Figure 5.5: Regularization effects using `min_samples_leaf`

5.2.1 Computational Complexity

To predict with a Decision Tree the algorithm follows a path from the root to a leaf. Because most trained trees are roughly balanced, this path crosses about

$$T_{\text{predict}} = O(\log_2 m), \quad (5.3)$$

nodes, where m is the number of training samples. Each node tests a single feature, so the prediction cost in equation 5.3 is unaffected by the total feature count and is therefore very fast.

Training is costlier. At every split the algorithm scans each candidate feature (or the subset limited by `max_features`) over all samples that reach the node. For n input features the total work across all levels is

$$T_{\text{train}} = O(n \times m \log m), \quad (5.4)$$

because the tree has roughly $\log_2 m$ layers and each layer processes a shrinking fraction of the m samples.

When the dataset contains only a few thousand samples, enabling `presort=True` in Scikit-Learn can shorten training time, but for larger datasets the presorting step turns into a bottleneck and the standard (unsorted) approach is quicker.

5.2.2 Entropy verse Gini Impurity

Decision Trees in sklearn use the Gini impurity as the default criterion for node purity, entropy can also be utilized by setting the criterion hyperparameter to “entropy”. Originally a thermodynamic concept representing molecular disorder, entropy reaches zero when molecules are in a complete state of rest and order. This concept was later adopted in information theory, introduced by Shannon, to quantify the average informational content in messages entropy is zero when all messages are identical. In the realm of Machine Learning, entropy serves as a measure of impurity: it is zero in a dataset if all its elements belong to a single class. The entropy for the i^{th} node is defined as

$$H_i = - \sum_{\substack{k=1 \\ p_{i,k} \neq 0}}^n p_{i,k} \log_2(p_{i,k}). \quad (5.5)$$

The choice between Gini impurity and entropy often results in negligible differences, producing similar trees. Gini impurity has a slight computational advantage and is thus the default choice. However, it tends to separate the most frequent class into a distinct branch, whereas entropy generally yields more balanced trees.

5.3 Decision Tree Regression

Decision Trees are not limited to classification tasks; they can also be adapted for regression. To illustrate, we use Scikit-Learn’s `DecisionTreeRegressor` to train a regression tree on a noisy linear dataset with `max_depth=2`. The structure of this tree is displayed in Figure 5.6. In contrast to the color scheme for classification, the change in color shades here relate to the predicted value of y .

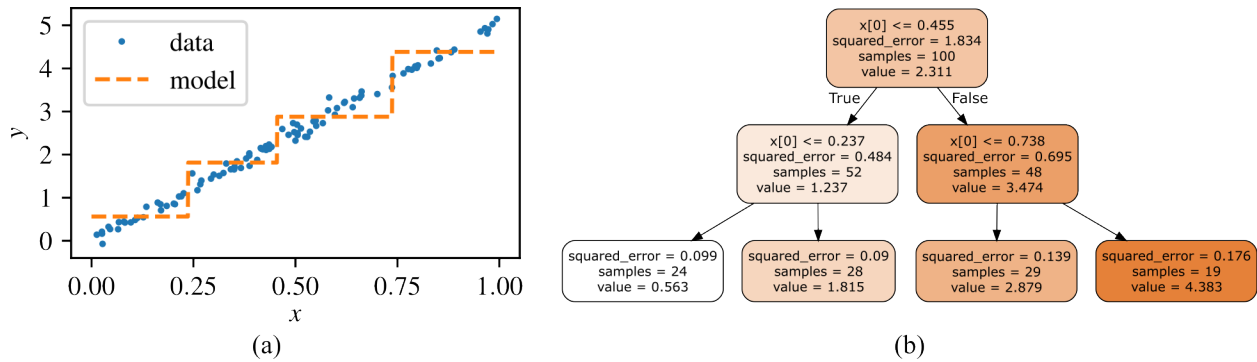


Figure 5.6: A regression model developed using a Decision Tree, showing the: (a) model superimposed over noisy data, and (b) the decision tree developed for the task.

The workings of a regression tree mirror those of a classification tree, except that each leaf holds a continuous prediction instead of a class label. To obtain a prediction for a sample with $x_1 = 0.6$, trace the path from the root to its leaf; that leaf outputs 0.1106, which is the average target value among the 110 training instances that reach it. For those instances the mean squared error is 0.0151.

The predictions from this model are visualized in Figure 5.7, with results shown for trees of depth 2 and 3. The deeper tree partitions the input space into more regions, with each region's predicted value being the average target value of the instances it encompasses. The model attempts to organize the regions such that the instances within each are as close as possible to their predicted value.

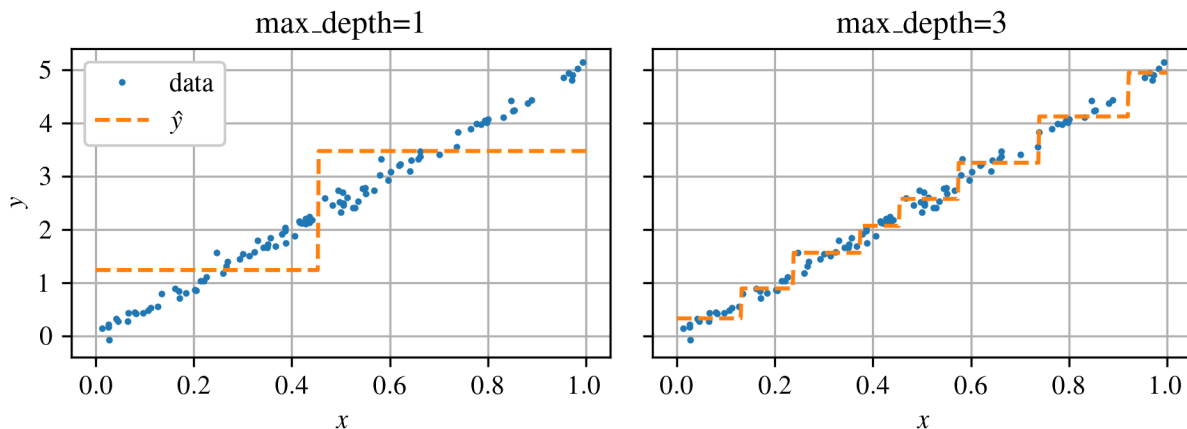


Figure 5.7: Comparison of predictions from two Decision Tree regression models with varying depths.

The CART algorithm for regression trees aims to minimize the MSE when splitting the training set, similar to how it minimizes impurity in classification tasks. The cost function minimized by the algorithm is represented as

$$J(k, t_k) = \frac{m_{\text{left}}}{m} \text{MSE}_{\text{left}} + \frac{m_{\text{right}}}{m} \text{MSE}_{\text{right}}. \quad (5.6)$$

Knowing that,

$$\text{MSE}_{\text{node}} = \sum_{i \in \text{node}} \left(\hat{y}_{\text{node}} - y^{(i)} \right)^2 \quad (5.7)$$

and

$$\hat{y}_{\text{node}} = \frac{1}{m_{\text{node}}} \sum_{i \in \text{node}} y^{(i)}. \quad (5.8)$$

Similar to classification, Decision Trees for regression can overfit if not properly regularized. Without regularization, the predictions, as depicted on the left of Figure 5.8, can fit the training data excessively. By setting `min_samples_leaf` to 15, a more generalized model is achieved, as shown on the right in the same figure.

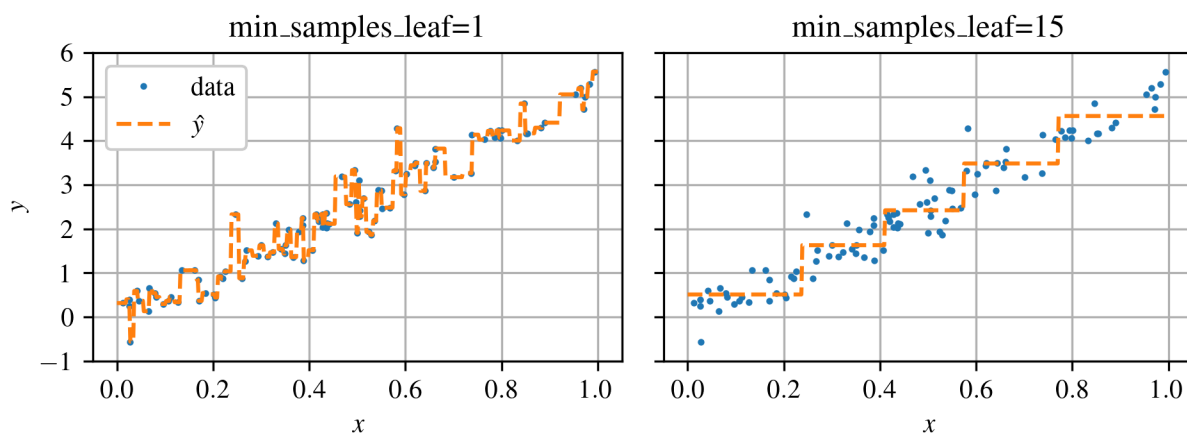


Figure 5.8: Impact of regularization through setting the minimum number of samples in a leaf on a Decision Tree regression model.

Example 5.2 Decision Tree Regression

This example uses Scikit-Learn's `DecisionTreeRegressor` to fit a non-linear relationship between input and output data. A tree of depth 4 is trained on noisy linear data and visualized using Graphviz to show how the regression model partitions the input space.

5.4 Instability

While Decision Trees offer simplicity, interpretability, versatility, and power, they come with certain drawbacks. A notable limitation is their preference for orthogonal decision boundaries, which makes them highly sensitive to the orientation of the data. For instance, Figure 5.9 illustrates how a simple linearly separable dataset can be perfectly split by a Decision Tree in its original alignment, whereas a 45° rotation results in a convoluted decision boundary. Despite perfect fits to their respective training sets, the rotated tree's model is likely to perform poorly on unseen data. Utilizing Principal Component Analysis (PCA) can often mitigate this issue by reorienting the data more suitably.

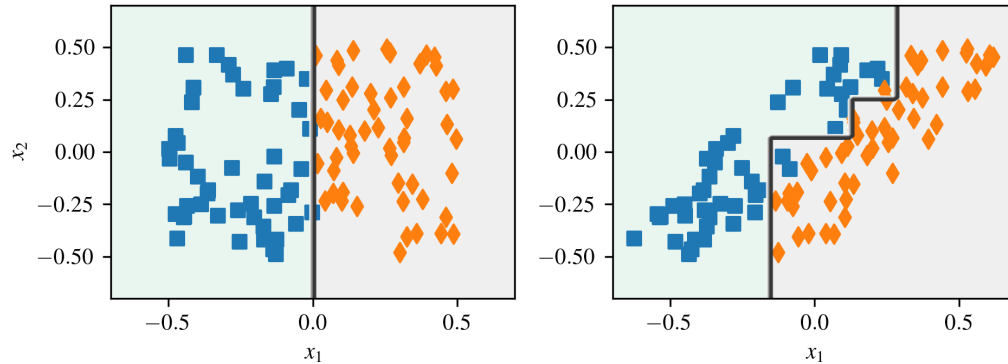


Figure 5.9: Sensitivity to training set rotation. Decision trees create a single clean split on the original data (left), but after a 45° rotation they must build a jagged, multi-step boundary; illustrating their sensitivity to feature orientation.

More broadly, Decision Trees are sensitive to small variations in training data. For instance, changing the seed of the random number generator can lead to a substantially different model, as shown in Figure 5.10. This variability is partly due to the stochastic nature of the Scikit-Learn’s training algorithms; different models may result from the same data unless the `random_state` hyperparameter is fixed.

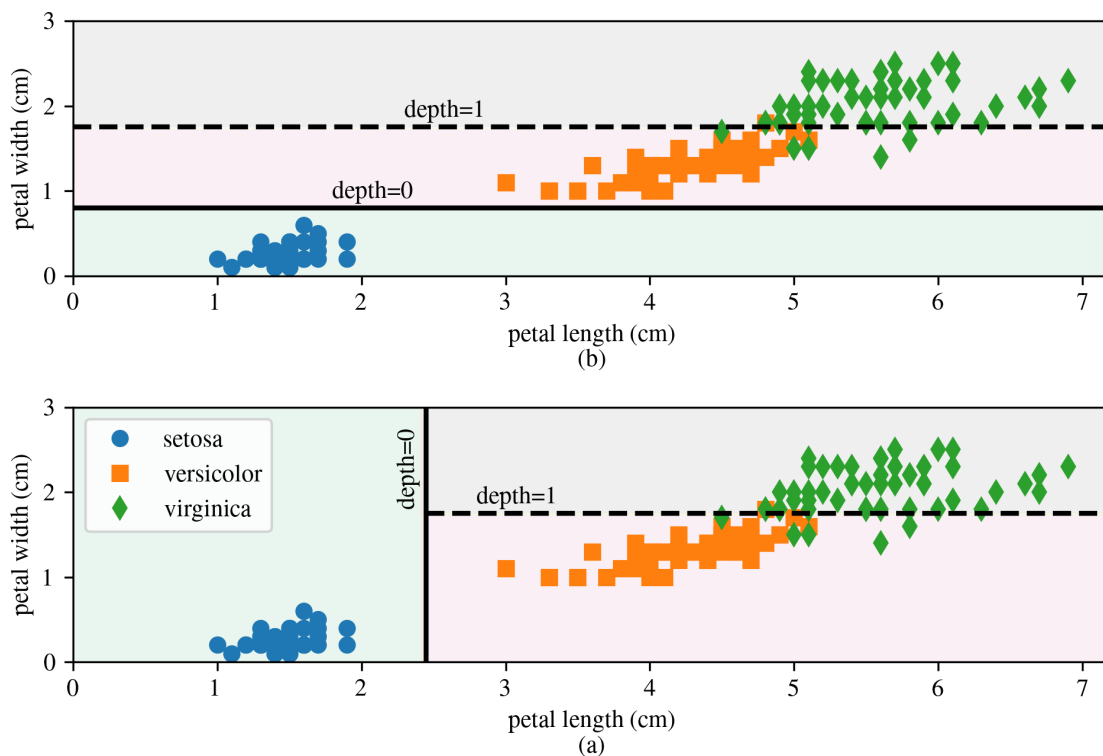


Figure 5.10: Decision Tree Sensitivity to initial conditions, showing: (a) random number generator seeded with “1”, and: (b) random number generator seeded with “2”.^a

5.5 Random Forest

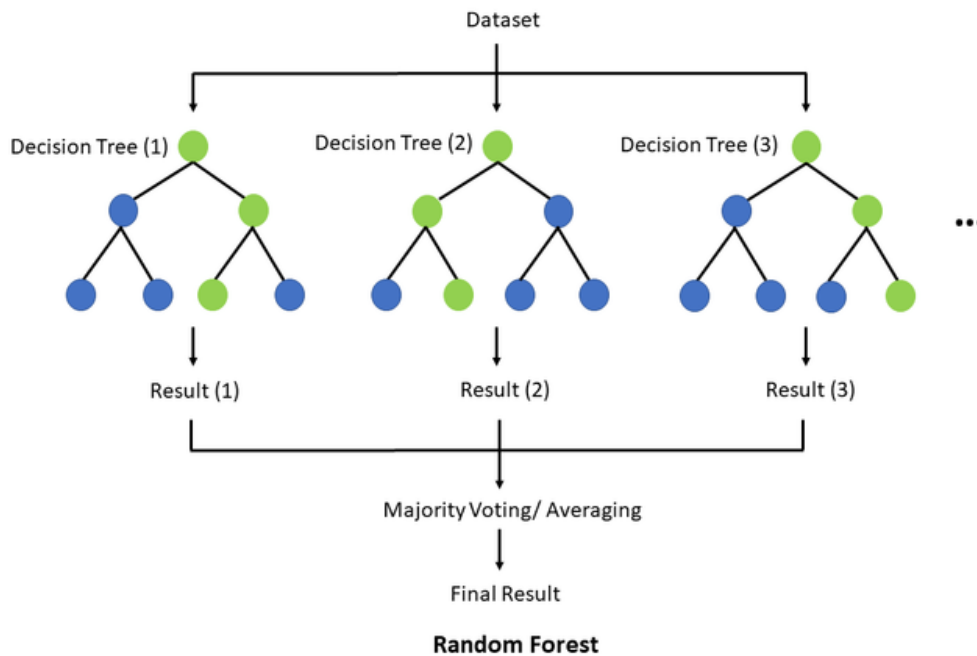


Figure 5.11: Random forest with majority voting. ^a

Random Forests reduce the high variance of individual decision trees by building many trees on bootstrap samples of the training data and then averaging their outputs. Each tree is grown with a random subset of input features, so the ensemble decorrelates the trees and improves generalisation. For classification problems the forest predicts the class that receives the majority vote, while for regression it returns the mean of the trees' numeric predictions. This bagging strategy limits overfitting and usually yields better performance than a single tree, although it often falls short of the accuracy achieved by gradient-boosted ensembles. Model quality still depends on the data's size, noise level, and feature interactions, and on hyperparameters such as the number of trees, the maximum depth, and the number of features considered at each split.

^a“2” is the author's preferred random number seed.

^aTseKiChun, CC BY-SA 4.0 <<https://creativecommons.org/licenses/by-sa/4.0/>>, via Wikimedia Commons

5.6 Examples

Example 5.1

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Example 5.1 Decision Tree Classifier
5
6  @author: Austin R.J. Downey
7  """
8
9  import IPython as IP
10 IP.get_ipython().run_line_magic('reset', '-sf')
11
12 import numpy as np
13 import matplotlib.pyplot as plt
14 import sklearn as sk
15 from sklearn.datasets import load_iris
16 from sklearn.tree import export_graphviz
17 import graphviz as graphviz
18
19 plt.close('all')
20
21
22 """ Load your data
23
24 # We will use the Iris data set. This dataset was created by biologist Ronald
25 # Fisher in his 1936 paper "The use of multiple measurements in taxonomic
26 # problems" as an example of linear discriminant analysis
27 iris = sk.datasets.load_iris()
28
29 # for simplicity, extract some of the data sets
30 X = iris['data'] # this contains the length of the petals and sepals
31 Y = iris['target'] # contains what type of flower it is
32 Y_names = iris['target_names'] # contains the name that aligns with the type of the
33 flower
34 feature_names = iris['feature_names'] # the names of the features
35
36 """ Build the model
37
38 # train the decision tree
39 tree_clf = sk.tree.DecisionTreeClassifier(max_depth=3)
40 X_petal = X[:,2:]
41 tree_clf.fit(X_petal, Y)
42
43 """ Visualize the decision tree
44
45 #create the export file for graphviz and export it. The file is exported as a
46 #.DOT file and can be viewed in an online viewer
47 https://dreampuf.github.io/GraphvizOnline/
48 export_graphviz(
49     tree_clf,
50     out_file="tree_clf.dot",
51     feature_names=feature_names[2:],
52     class_names=Y_names,
53     rounded=True,
54     filled=True
55 )
56
57 # We can load the file back in
58 s = graphviz.Source.from_file('tree_clf.dot')
59
60 # look at what is inside it. Also, just typing s in the console will display the image
61 print(s)
62
63 # export the image to a jpg
64 s.render('tree_clf', format='jpg',view=True)
65
66 """ Predict the class for any petal size

```

```
66
67     size = [[7, 2.5]]
68     print(tree_clf.predict_proba(size))
69     print(iris.target_names)
70
71
72     # plot the new data point over the Iris dataset
73     plt.figure()
74     plt.grid(True)
75     plt.scatter(X[Y==0,2],X[Y==0,3],marker='o',label=Y_names[0],zorder=2)
76     plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1],zorder=2)
77     plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2],zorder=2)
78     plt.scatter(size[0][0],size[0][1],s=300,marker='*',label='new data point',zorder=2)
79     plt.xlabel(feature_names[2])
80     plt.ylabel(feature_names[3])
81     plt.legend(framealpha=1)
82     plt.tight_layout()
83
84
85
86
87
```

Example 5.2

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Example 5.2 Decision Tree Regression
5
6  @author: Austin R.J. Downey
7  """
8
9  import IPython as IP
10 IP.get_ipython().run_line_magic('reset', '-sf')
11
12 import numpy as np
13 import matplotlib.pyplot as plt
14 import sklearn as sk
15 import graphviz as graphviz
16 from sklearn.tree import export_graphviz
17
18 plt.close('all')
19
20
21 """ Train and plot a decision tree regression model
22
23 # build the data
24 m = 200
25 X = np.random.rand(m, 1)
26 y = 5 * X
27 y = y + np.random.randn(m, 1) / 10
28
29 # train the model
30 tree_reg = sk.tree.DecisionTreeRegressor(max_depth=3)
31 tree_reg.fit(X, y)
32
33 x_model = np.linspace(0, 1, 100).reshape(-1, 1)
34 y_model = tree_reg.predict(x_model)
35
36 plt.figure()
37 plt.plot(X, y, ".", label='data')
38 plt.plot(x_model, y_model, "-", label="model")
39 plt.xlabel("$x$")
40 plt.ylabel("$y$")
41 plt.legend()
42
43 """ Plot the regression tree
44
45 # create the export file for graphviz and export it. The file is exported as a
46 # .DOT file and can be viewed in an online viewer
47 https://dreampuf.github.io/GraphvizOnline/
48 export_graphviz(
49     tree_reg,
50     out_file="tree_reg",
51     rounded=True,
52     filled=True
53 )
54
55 # We can load the file back in
56 s = graphviz.Source.from_file('tree_reg')
57 s.render('tree_reg', format='jpg', view=True)
58
59
60
61

```

6 Support Vector Machines

The core idea of Support Vector Machines becomes clear when viewed geometrically. Figure 6.1 shows a slice of the iris dataset containing two linearly separable species. In the left panel two conventional linear classifiers are drawn; although both label every training point correctly, their separating lines sit very close to several observations, so even a small perturbation could lead to misclassifications on new data. The right panel plots the boundary produced by an SVM. The solid line still splits the classes cleanly, but it is positioned to maximise the distance to the nearest points, leaving the widest possible “street” (bounded by the dashed parallels). The samples that lie on these margins are the *support vectors*. This strategy, called large margin classification, usually delivers models that generalise better than those that merely separate the training set.

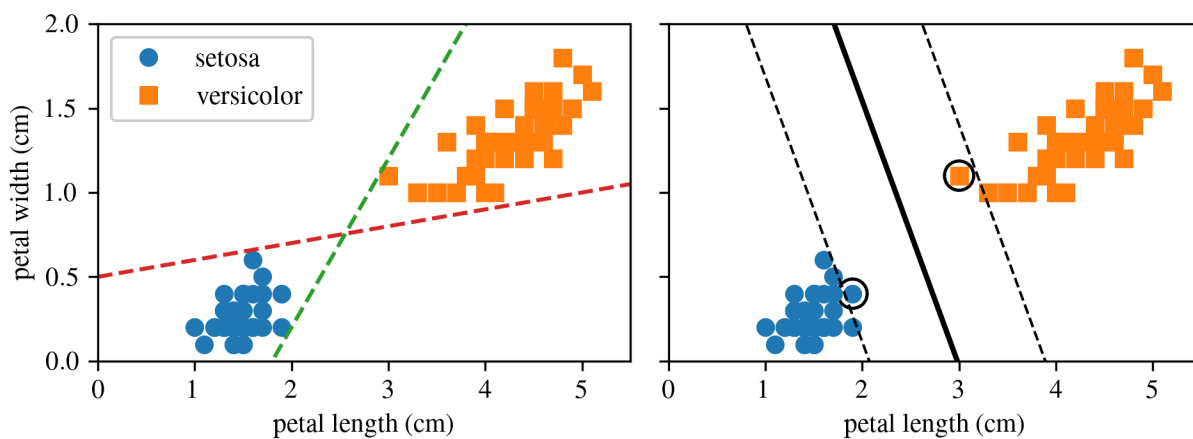
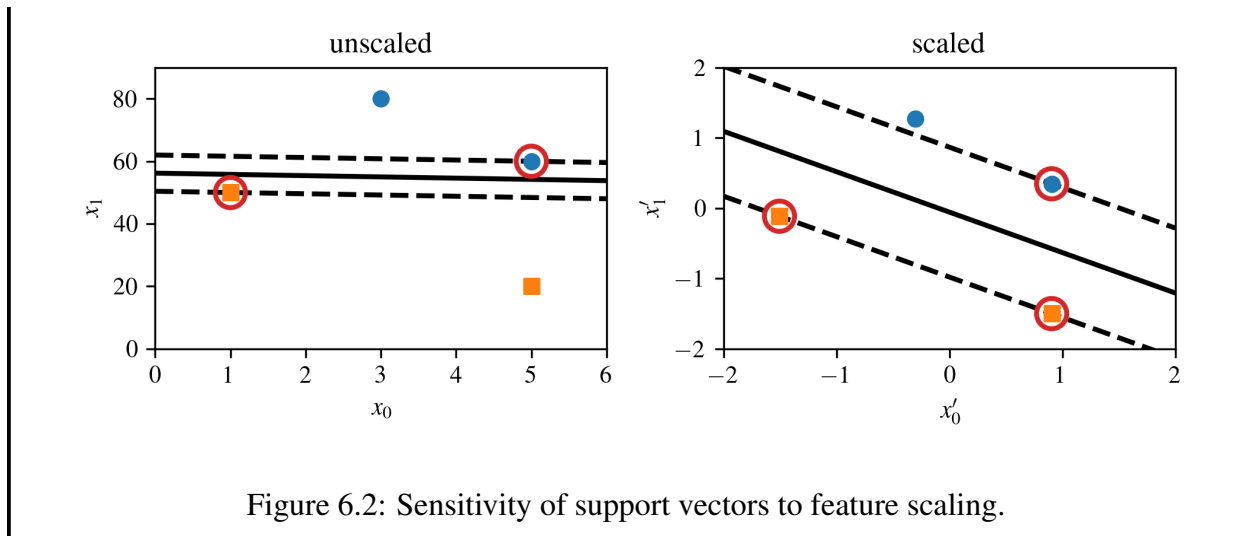


Figure 6.1: Large margin classification.

Observe that the addition of further training instances outside the “street” does not influence the decision boundary; it is entirely shaped by the instances situated on the boundary’s edge. These pivotal instances are termed support vectors and are highlighted with circles in Figure 6.1.

NOTE

The sensitivity of SVMs to feature scales is evident in Figure 6.2. In the left plot, the vertical dimension greatly outweighs the horizontal dimension, resulting in a nearly horizontal “street.” However, after applying feature scaling such as using Scikit-Learn’s `StandardScaler` the decision boundary becomes more appropriate, as illustrated in the right plot.



6.1 Linear SVM Classification

Hard margin classification demands that all instances be correctly classified without any margin violations. This strict approach faces two significant challenges:

- It is only feasible when the data is linearly separable.
- It is highly sensitive to outliers.

Figure 6.3 illustrates these challenges using the iris dataset with an added outlier. On the left, achieving a hard margin is impossible due to the outlier. On the right, although a decision boundary is found, it deviates substantially from the optimal boundary shown in Figure 6.1 and is less likely to perform well on new data.

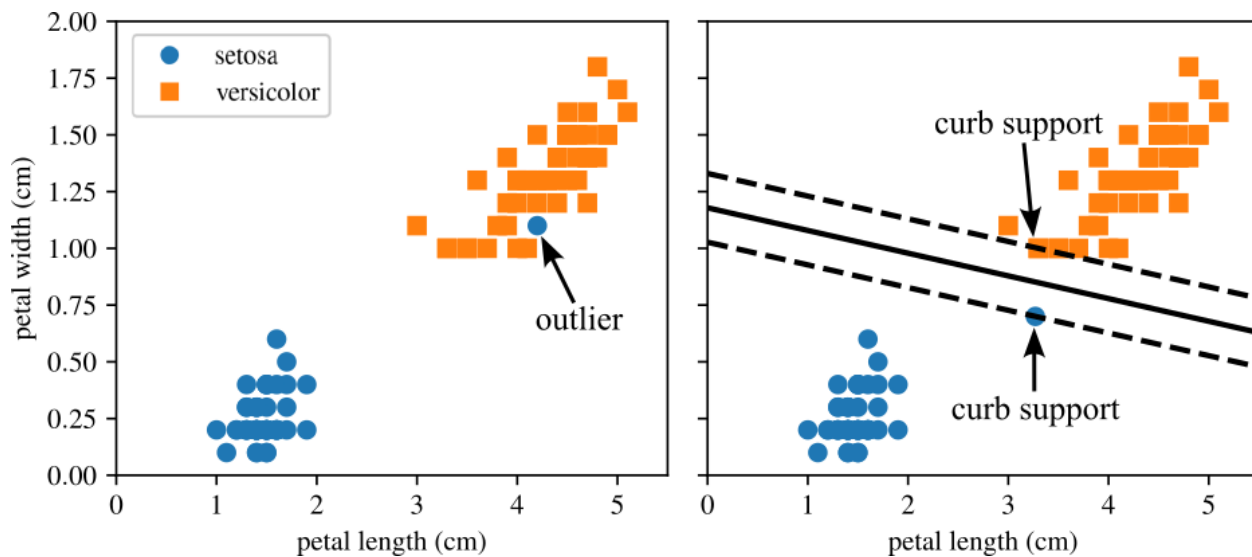


Figure 6.3: Support vector machines showing (left) an un-separable case, and (right) a separable case with two data points supporting the curbs of the support vector machine.

To mitigate the limitations of hard margin classification, a more adaptable model, known as soft margin classification, is often employed. The goal here is to achieve an optimal balance between maximizing the margin width and minimizing margin violations, where instances might fall into the margin or on the incorrect side.

Scikit-Learn's SVM implementations facilitate this balance through the hyperparameter C . A smaller value of C results in a wider margin but allows more margin violations, which is beneficial for model flexibility. Conversely, a larger C value tightens the margin, reducing margin violations but at the risk of a less flexible model. Figure 6.4 demonstrates this trade-off: the left plot with a low C value

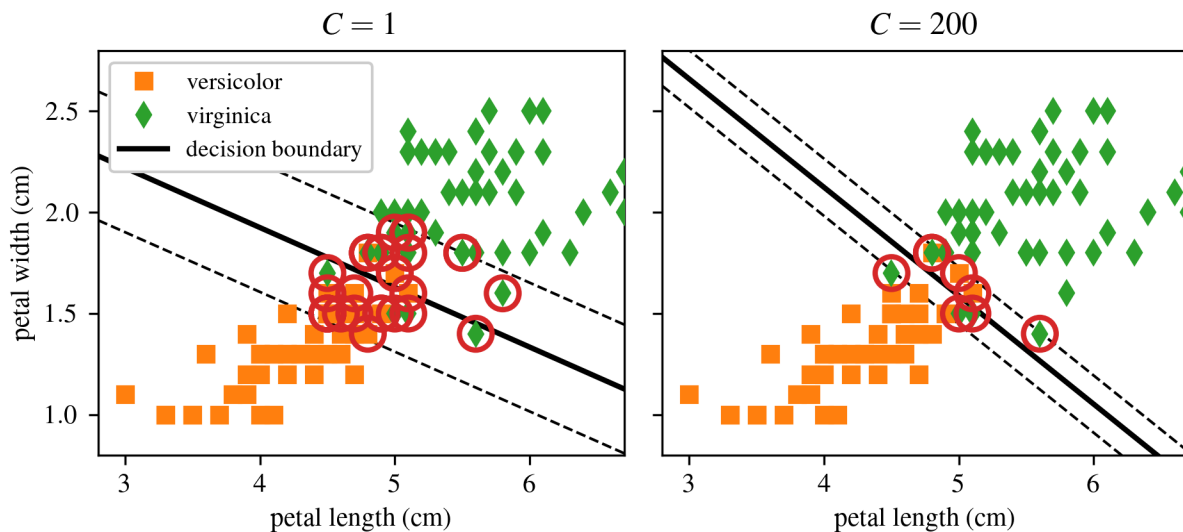


Figure 6.4: SVM margin sizes for different C values.

NOTE

Overfitting in an SVM model can often be addressed by reducing the C value, which increases regularization.

In earlier chapters we placed every model parameter in a single vector θ : the first entry θ_0 acted as the bias, while $\theta_1, \dots, \theta_n$ were the feature weights, and each input was augmented with a constant bias feature $x_0 = 1$. In this chapter we adopt the notation most common for SVMs. The bias is written as b , the weight vector as \mathbf{w} , and no extra bias feature is appended to the input vectors.

6.1.1 Decision Function and Predictions

A linear SVM classifier determines the class of a new instance x by calculating the decision function

$$\mathbf{w}^T \mathbf{x} + b = w_1 x_1 + \dots + w_n x_n + b. \quad (6.1)$$

If the outcome is positive, the predicted class \hat{y} is the positive class (1); otherwise, it is the negative class (0). This is written as

$$\hat{y} = \begin{cases} 0 & \text{if } \mathbf{w}^T \mathbf{x} + b < 0, \\ 1 & \text{if } \mathbf{w}^T \mathbf{x} + b \geq 0. \end{cases} \quad (6.2)$$

Figure 6.5 plots the decision function for a model with two features, so the surface is a plane in \mathbb{R}^2 . The thick solid line marks the decision boundary where the function equals zero. The dashed lines show the loci where the function equals 1 and -1 ; they run parallel to the boundary and sit at equal distance from it, outlining the margin. Training a linear SVM adjusts \mathbf{w} and b to make this margin as wide as possible while either prohibiting margin violations (hard margin) or keeping them small through a penalty term (soft margin).

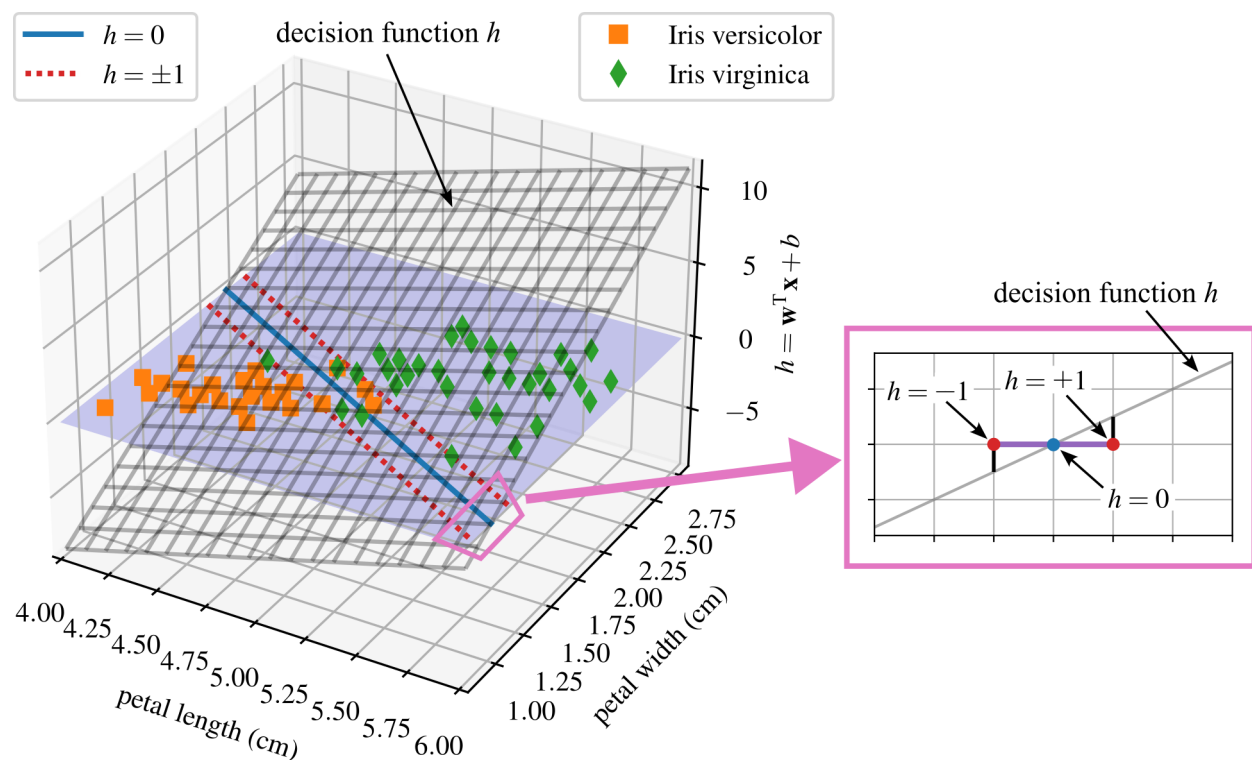


Figure 6.5: Decision function for the Iris Dataset showing how the decision function h cuts through the feature space.

6.1.2 Training Objective

The slope of the decision function corresponds to the norm of the weight vector, $\|\mathbf{w}\|$. Halving this slope causes the decision boundary margins, where the decision function equals ± 1 , to double in distance from the decision boundary. Effectively, reducing the norm of \mathbf{w} by half doubles the margin. This geometric interpretation is perhaps simpler to visualize in two dimensions, as shown in Figure 6.6. Thus, minimizing $\|\mathbf{w}\|$ maximizes the margin.

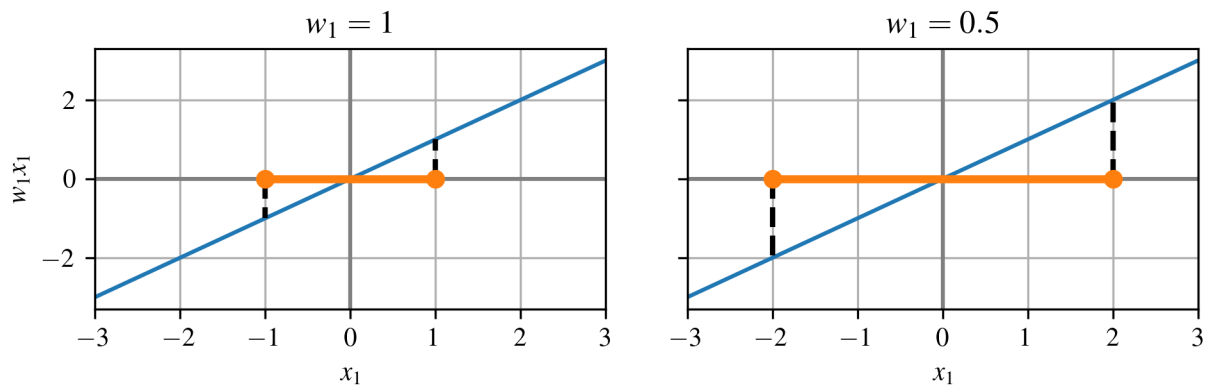


Figure 6.6: The margin is dependent on the value of the weight vector where a smaller weight vector results in a larger margin and vice versa.

To achieve a large margin while enforcing that no data points fall within the margin (hard margin), we ensure the decision function exceeds $+1$ for all positive training instances and is less than -1 for all negative instances. Let $t^{(i)}$ equal -1 for negative instances ($y^{(i)} = 0$) and $+1$ for positive ones ($y^{(i)} = 1$). The constraints then require

$$t^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \quad (6.3)$$

for all training instances. This forms the basis of the hard margin linear SVM classifier optimization problem:

$$\begin{aligned} & \underset{\mathbf{w}, b}{\text{minimize}} && \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ & \text{subject to} && t^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 \text{ for } i = 1, 2, \dots, m \end{aligned} \quad (6.4)$$

NOTE

The objective function minimized is $\frac{1}{2} \mathbf{w}^T \mathbf{w}$, equivalent to $\frac{1}{2} \|\mathbf{w}\|^2$. This formulation is chosen over minimizing $\|\mathbf{w}\|$ directly because $\frac{1}{2} \|\mathbf{w}\|^2$ offers a straightforward derivative, simply \mathbf{w} , facilitating gradient calculations. In contrast, $\|\mathbf{w}\|$ lacks differentiability at $\mathbf{w} = 0$, posing challenges for optimization algorithms, which typically require smooth, differentiable functions to ensure effective optimization.

To formulate the soft margin objective, it is necessary to introduce a slack variable $\zeta^{(i)} \geq 0$ for each instance. This variable, $\zeta^{(i)}$, quantifies the permissible margin violation for the i^{th} instance. Consequently, we face dual objectives: minimizing the slack variables to reduce margin violations and minimizing $\frac{1}{2}\mathbf{w}^T\mathbf{w}$ to maximize the margin. The hyperparameter C plays a crucial role here, enabling a balance between these competing objectives. The introduction of C transforms our task into a constrained optimization problem.

$$\begin{aligned} & \underset{\mathbf{w}, b, \zeta}{\text{minimize}} && \frac{1}{2}\mathbf{w}^T\mathbf{w} + C \sum_{i=1}^m \zeta^{(i)} \\ & \text{subject to} && t^{(i)}(\mathbf{w}^T\mathbf{x}^{(i)} + b) \geq 1 - \zeta^{(i)} \text{ and } \zeta^{(i)} \geq 0 \text{ for } i = 1, 2, \dots, m \end{aligned} \quad (6.5)$$

6.1.3 Quadratic Programming

Both hard margin and soft margin problems are examples of convex quadratic optimization problems with linear constraints, commonly referred to as Quadratic Programming (QP) problems. QP involves solving optimization problems where the objective is a quadratic function and the constraints are linear. This form of programming, established in the 1940s, predates and is distinct from “computer programming,” and is sometimes more descriptively termed “quadratic optimization” to avoid confusion.

A variety of techniques available through off-the-shelf solvers can address these QP problems, though they extend beyond the scope of this text. The general formulation of a QP problem is as follows:

$$\begin{aligned} & \underset{\mathbf{p}}{\text{minimize}} && \frac{1}{2}\mathbf{p}^T\mathbf{H}\mathbf{p} + \mathbf{f}^T\mathbf{p} \\ & \text{subject to} && \mathbf{A}\mathbf{p} \leq \mathbf{b} \end{aligned} \quad (6.6)$$

Here, \mathbf{p} is an n_p -dimensional vector (where n_p is the number of parameters), \mathbf{H} is an $n_p \times n_p$ matrix, \mathbf{f} is an n_p -dimensional vector, \mathbf{A} is an $n_c \times n_p$ matrix (with n_c being the number of constraints), and \mathbf{b} is an n_c -dimensional vector.

Equation 6.6 defines a standard quadratic program with constraints of the form $\mathbf{A}\mathbf{p} \leq \mathbf{b}$. For training a hard margin linear SVM, this setup can be used by choosing parameters that encode the SVM objective and constraints. The solution vector \mathbf{p} contains both the bias term and the feature weights. Using this knowledge, a standard QP solver can be applied directly to find the optimal SVM classifier.

Example 6.1 Support Vector Machine Classification

This example applies a linear Support Vector Machine (SVM) classifier to distinguish Iris-Virginica flowers based on petal length and width. The decision boundary is derived after scaling the data, and margins are visualized. Misclassified instances within the margin are identified and marked. The model’s performance is assessed using a confusion matrix and F1 score.

6.2 Nonlinear SVM Classification

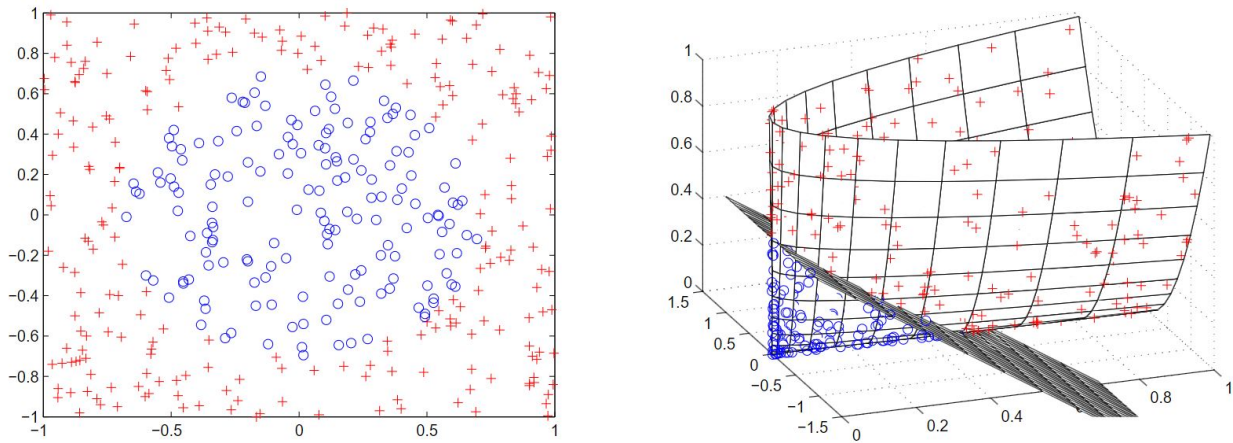


Figure 6.7: Nonlinear SVM example illustration ^a

While linear SVM classifiers are quite effective and perform exceptionally well in various scenarios, many datasets are far from being linearly separable. One strategy to address non-linear datasets is to introduce additional features, such as polynomial feature. Adding features can sometimes transform the dataset into one that is linearly separable. A representation of this technique is shown in 6.7.

A simple example of converting non-linearly separable variables is shown in figure 6.8 where the left plot displays a simple dataset with a single feature x_1 . Clearly, this dataset is not linearly separable. However, by adding another feature $x_2 = (x_1)^2$, the dataset becomes perfectly linearly separable in two dimensions.

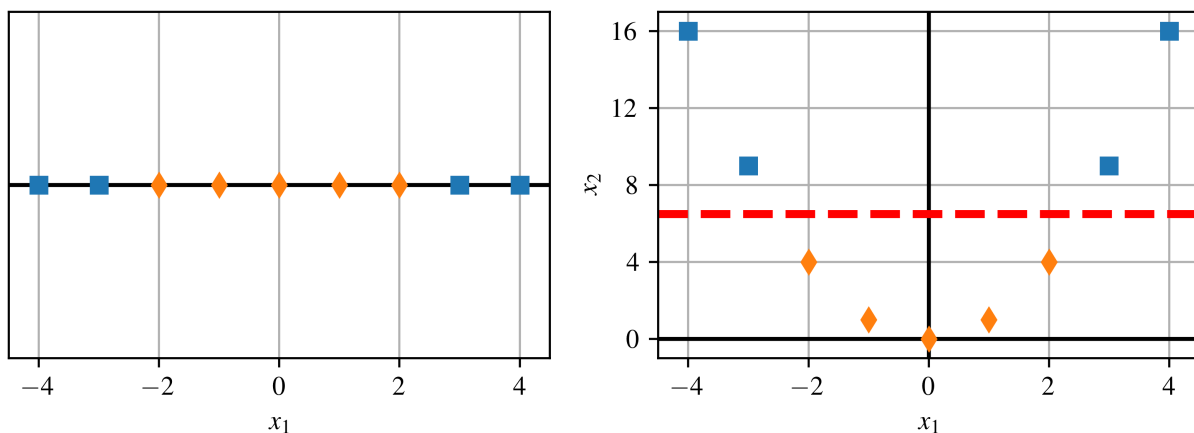


Figure 6.8: Illustration of SVM in higher dimensions

^aMachine Learner, CC BY-SA 4.0 <<https://creativecommons.org/licenses/by-sa/4.0/>>, via Wikimedia Commons

You can implement this idea in Scikit-Learn by creating a pipeline that applies a Polynomial Features transformer (introduced in the Regression Chapter), followed by a `StandardScaler` and a `LinearSVC`. The approach works nicely on the moons dataset, a toy binary-classification problem in which the samples trace two interleaving half-circles, as illustrated in Figure 6.9. You can generate this dataset with the function `make_moons()`.

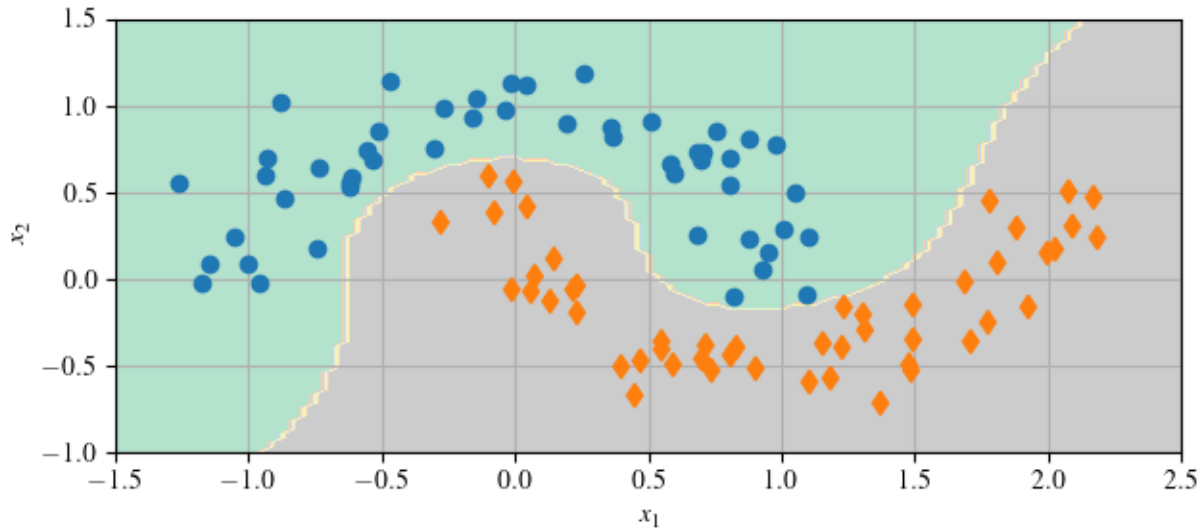


Figure 6.9: Demonstration of a SVM classifier with polynomial features.

Example 6.2 Nonlinear Classification

This example demonstrates nonlinear classification leveraging `LinearSVC` in a pipeline with `preprocessing.PolynomialFeatures` using the `make_moons` dataset. A polynomial feature transformation is combined with a linear SVM to classify the data, and the resulting decision boundaries are visualized.

6.2.1 Kernel trick

While adding polynomial features is straightforward and enhances the performance of various Machine Learning algorithms (not limited to SVMs), it presents limitations. Specifically, lower-degree polynomials may not adequately handle complex datasets, and higher-degree polynomials significantly increase the feature count, slowing down the model.

SVMs offer a unique solution through a remarkable mathematical technique known as the kernel trick, which allows for the benefits of high-degree polynomial features without actually expanding the feature space, thereby avoiding a rapid increase in computation. This kernel trick is incorporated within the `SVC` class.

In the dual form of an SVM the explicit dot product of two input vectors

$$\mathbf{x}^T \mathbf{z} \quad (6.7)$$

is replaced by a kernel function

$$k(\mathbf{x}, \mathbf{z}), \quad (6.8)$$

where \mathbf{x} and \mathbf{z} are n -dimensional feature vectors representing any two data points in the training set. This substitution lets the algorithm construct a linear separator in a (possibly infinite-dimensional) feature space while all calculations still occur in the original input coordinates.

Figure 6.10 shows configured SVM classifiers using a 3rd and 3th degree polynomial kernel; on the left and right respectively. Adjusting the polynomial degree can help manage model fit: reducing the degree may prevent overfitting, whereas increasing it may be necessary for underfitting scenarios. The ‘coef0’ hyperparameter is crucial as it determines the influence of high versus low-degree polynomials in the model.

NOTE

A typical method for determining optimal hyperparameter settings involves utilizing grid search techniques. Starting with a broad, coarse grid search to quickly narrow down potential candidates, followed by a more detailed, finer grid search centered on these promising values often yields faster results. Additionally, understanding the function and influence of each hyperparameter aids in efficiently targeting the most relevant areas of the hyperparameter space.

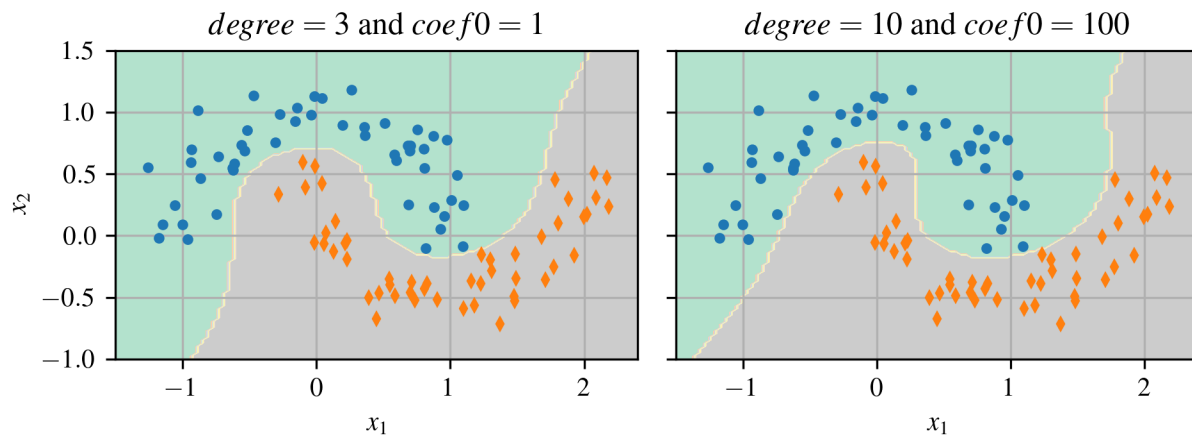


Figure 6.10: SVM polynomial kernel

6.2.2 Standard Kernels

Three kernels cover the vast majority of practical cases.

- The polynomial kernel captures interactions between input features up to a chosen degree d :

$$k_{\text{poly}}(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + c)^d \quad (6.9)$$

where $c \geq 0$ is a constant offset. It is effective when domain knowledge suggests that a low-order combination of variables explains the target. Keep d modest (typically $d \leq 5$) and standardise inputs to avoid exploding feature dimensions and overfitting.

- The radial basis function (RBF) kernel builds smooth, highly flexible decision boundaries:

$$k_{\text{rbf}}(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|^2) \quad (6.10)$$

with width parameter $\gamma > 0$. A large γ makes the surface too flat (underfitting), while a small γ lets every point carve its own pocket (overfitting). Tune C and γ jointly, typically on a logarithmic grid, after x-score standardising the features.

- The sigmoid kernel adds an S-shaped non-linearity to the inner product:

$$k_{\text{sig}}(\mathbf{x}, \mathbf{z}) = \tanh(\kappa \mathbf{x}^T \mathbf{z} + \theta) \quad (6.11)$$

where κ controls the slope and θ the offset. Although useful for some sparse or text data, this kernel is not always positive-semidefinite, so ensure your software handles the resulting optimisation safely.

Example 6.3 Polynomial Kernel Trick

This example uses the kernel trick to enable an SVM to classify non-linearly separable data using SVC (not LinearSVC). A third-degree polynomial kernel is applied to the moons dataset using Scikit-Learn's Pipeline and SVC tools, producing a curved decision boundary that cleanly separates the classes.

NOTE

Begin with the RBF kernel as a strong default, explore polynomial kernels when you expect specific interaction orders, and treat the sigmoid option as experimental unless you have evidence it helps.

6.3 Computational Complexity

Scikit-Learn provides two main SVM classifiers that trade accuracy for speed in different ways. LinearSVC is optimized for large, high-dimensional data when a linear decision boundary is adequate; SVC sacrifices runtime for the expressive power of kernels and thus handles complex, nonlinear patterns. Table 2 highlights how these priorities translate into computational costs and practical constraints and compares them to SGDClassifier for reference.

Table 2: Key characteristics of three Scikit-Learn SVM classifier implementations.

True class	time complexity	out-of-core support	scaling required	kernel trick
LinearSVC	$O(m \times n)$	no	yes	no
SVC	$O(m^2 \times n)$ to $O(m^3 \times n)$	no	yes	yes
SGDClassifier	$O(m \times n)$	yes	yes	no

LinearSVC builds on the liblinear solver and is limited to *linear* decision boundaries. Because it does not apply the kernel trick, its training cost grows almost linearly with both the number

of samples m and features n , i.e. $O(mn)$. Convergence is controlled by the tolerance parameter `tol` (denoted ε in the literature); the default value is usually sufficient, but smaller tolerances can be specified when higher accuracy is critical.

`SVC`, in contrast, relies on `libsvm` and *does* support kernel functions. Its computational burden is markedly heavier, between $O(m^2n)$ and $O(m^3n)$ in practice, so training becomes prohibitive once the dataset reaches the hundreds-of-thousands range. Nevertheless, `SVC` excels on smaller or medium-sized problems that demand nonlinear decision surfaces. Runtime also scales with the average count of non-zero features per instance, meaning sparse high-dimensional inputs remain tractable.

6.4 SVM Regression

Support Vector Regression (SVR) adapts the SVM idea to regression by surrounding the prediction curve with a tube of width ε ; points outside this tube incur a penalty and the optimizer keeps the tube as flat as possible while allowing a limited number of violations governed by the regularization constant C . Conceptually this reverses the classification goal: instead of widening a margin to separate two classes, SVR encourages the data to lie inside the margin, making ε the principal knob that controls how loose the fit may be for both linear and kernel-based models.

6.4.1 Linear SVR

The presence of additional training instances within the margin does not influence the predictions of the model, rendering it ε -insensitive. For linear SVM Regression tasks, the `LinearSVR` class from Scikit-Learn can be utilized. For instance, Figure 6.11 demonstrates two linear SVM Regression models trained on random linear data; one features a wide margin ($\varepsilon = 1.5$), and the other a narrower margin ($\varepsilon = 0.5$).

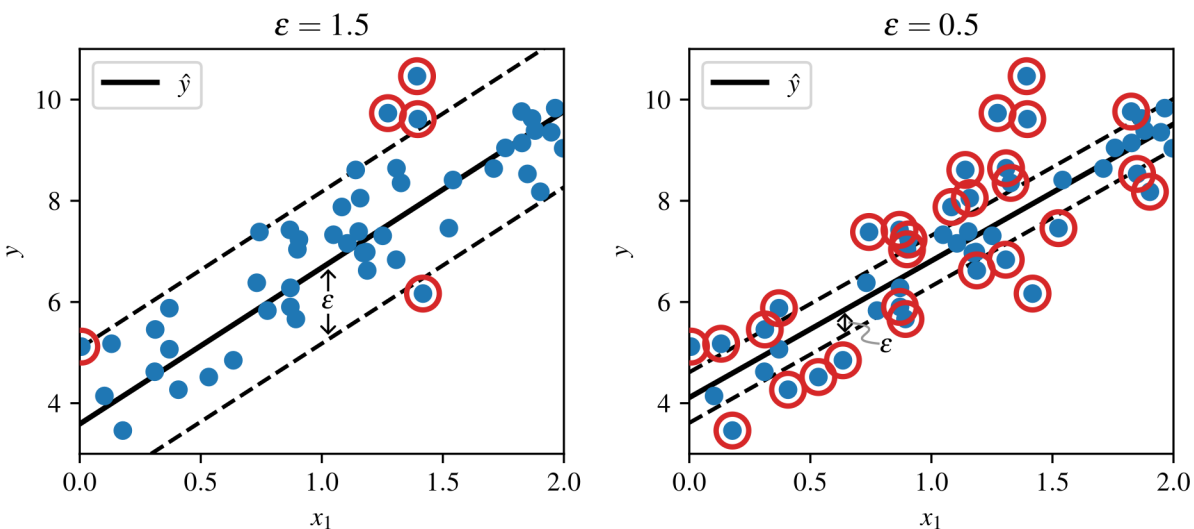


Figure 6.11: SVM Regression models with different ε values.

For data that follow an approximately linear trend, the `LinearSVR` class in Scikit-Learn solves the primal problem directly. It scales in $O(mn)$ time with m samples and n features, making it a practical choice for large data sets where a simple linear fit is adequate.

6.4.2 Kernel SVR

When the relationship is non-linear, the SVR class employs the kernel trick:

$$f(x) = \sum_{i=1}^m (\alpha_i - \alpha_i^*) k(x_i, x) + b, \quad (6.12)$$

where $k(\cdot, \cdot)$ is typically RBF or polynomial. Figure 6.12 shows a 2nd-degree polynomial kernel capturing quadratic structure under various regularisation levels.

The Scikit-learn SVR class, supporting the kernel trick and acting as the regression counterpart to the SVC class, performs well with small to medium-sized datasets but slows considerably as dataset size increases. In contrast, the LinearSVR class, akin to the LinearSVC class, scales linearly with the size of the training set.

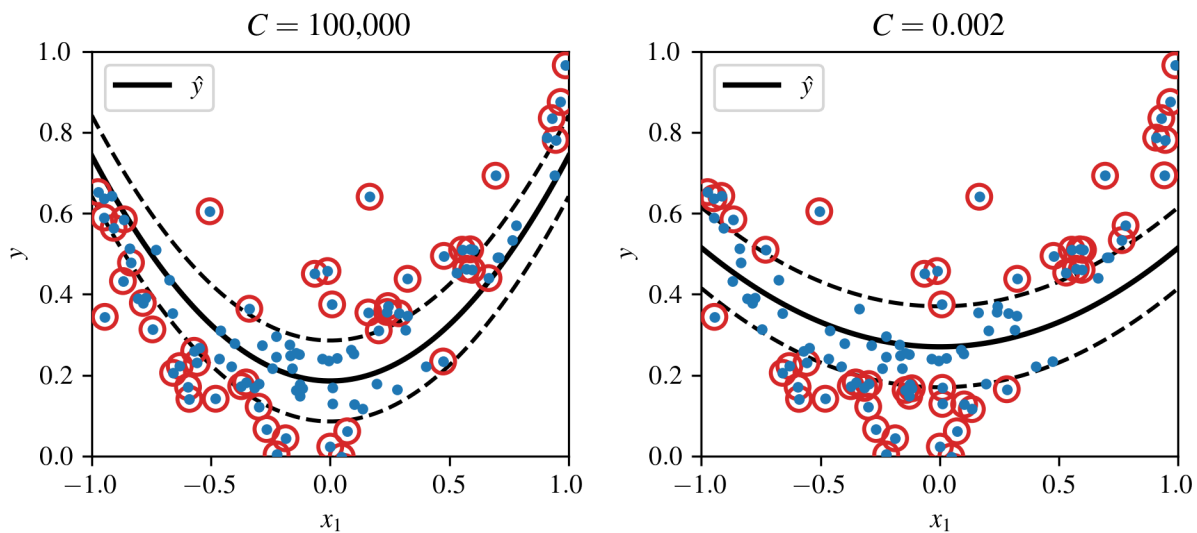


Figure 6.12: SVM regression with a 2nd-degree polynomial kernel, showcasing different regularization levels.

Example 6.4 SVM Polynomial Regression

This example fits a non-linear regression curve to noisy quadratic data using Support Vector Regression (SVR) with a polynomial kernel. It highlights how SVR models can handle non-linear relationships by introducing a margin of tolerance.

6.5 Examples

Example 6.1

```

1  """
2  Example 6.1 Support Vector Machine Classification
3  @author: Austin R.J. Downey
4  """
5
6  import IPython as IP
7  IP.get_ipython().run_line_magic('reset', '-sf')
8
9  import numpy as np
10 import matplotlib.pyplot as plt
11 import sklearn as sk
12 from sklearn import datasets
13 from sklearn import svm
14 from sklearn import pipeline
15
16 cc = plt.rcParams['axes.prop_cycle'].by_key()['color']
17 plt.close('all')
18
19
20 # This code will build and train a support vector machine classifier with soft
21 # for the iris flower data set.
22
23 %% Load your data
24
25 # We will use the Iris data set. This dataset was created by biologist Ronald
26 # Fisher in his 1936 paper "The use of multiple measurements in taxonomic
27 # problems" as an example of linear discriminant analysis
28 iris = sk.datasets.load_iris()
29
30 # for simplicity, extract some of the data sets
31 X = iris['data'] # this contains the length of the petals and sepals
32 Y = iris['target'] # contains what type of flower it is
33 Y_names = iris['target_names'] # contains the name that aligns with the type of the
34 # flower
35 feature_names = iris['feature_names'] # the names of the features
36
37 # plot the Sepal data
38 plt.figure(figsize=(6.5,3))
39 plt.subplot(121)
40 plt.grid(True)
41 plt.scatter(X[Y==0,0],X[Y==0,1],marker='o',zorder=10)
42 plt.scatter(X[Y==1,0],X[Y==1,1],marker='s',zorder=10)
43 plt.scatter(X[Y==2,0],X[Y==2,1],marker='d',zorder=10)
44 plt.xlabel(feature_names[0])
45 plt.ylabel(feature_names[1])
46
47 plt.subplot(122)
48 plt.grid(True)
49 plt.scatter(X[Y==0,2],X[Y==0,3],marker='o',label=Y_names[0],zorder=10)
50 plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1],zorder=10)
51 plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2],zorder=10)
52 plt.xlabel(feature_names[2])
53 plt.ylabel(feature_names[3])
54 plt.legend(framealpha=1)
55 plt.tight_layout()
56
57 %% Extract just the petal space of the code
58 X_petal = X[50:, (2, 3)] # petal length, petal width
59 y_petal = Y[50:] == 2
60
61 %% Build and train the SVM classifier
62
63 # build handles to regularize the model data and a Linear Support Vector Classification.
64 scaler = sk.preprocessing.StandardScaler()
65 svm_clf = sk.svm.LinearSVC(C=1000000,max_iter=10000)
66

```

```

67 # build the model pipeline of regularization and a Linear Support Vector Classification.
68 scaled_svm_clf = sk.pipeline.Pipeline([
69     ("scaler", scaler),
70     ("linear_svc", svm_clf),
71 ])
72
73 # train the data
74 scaled_svm_clf.fit(X_petal, y_petal)
75
76
77 %% Build and plot the decision boundary along with the curbs
78
79 # Convert to unscaled parameters as the SVM is solved in a scaled space.
80 w = svm_clf.coef_[0] / scaler.scale_
81 b = svm_clf.decision_function([-scaler.mean_ / scaler.scale_])
82
83 # At the decision boundary,  $w_0x_0 + w_1x_1 + b = 0$ 
84 # =>  $x_1 = -w_0/w_1 * x_0 - b/w_1$ 
85 x0 = np.linspace(4, 5.9, 200)
86 decision_boundary = -w[0]/w[1] * x0 - b/w[1]
87
88 margin = 1/w[1]
89 curbs_up = decision_boundary + margin
90 curbs_down = decision_boundary - margin
91
92 %% Plot the data and the classifier
93
94 plt.figure()
95 plt.grid(True)
96 plt.scatter(X[Y==1,2],X[Y==1,3],marker='s',label=Y_names[1],zorder=10)
97 plt.scatter(X[Y==2,2],X[Y==2,3],marker='d',label=Y_names[2],zorder=10)
98 plt.xlabel(feature_names[2])
99 plt.ylabel(feature_names[3])
100 plt.legend(framealpha=1)
101 plt.tight_layout()
102
103 # plot the decision boundy and margins
104 plt.plot(x0, decision_boundary, "k-", linewidth=2)
105 plt.plot(x0, curbs_up, "k--", linewidth=2)
106 plt.plot(x0, curbs_down, "k--", linewidth=2)
107
108
109 %% Find the misclassified instances and add a circle to mark them
110
111 # Find support vectors (LinearSVC does not do this automatically) and add them
112 # to the SVM handle
113 t = y_petal * 2 - 1 # convert 0 and 1 to -1 and 1
114 support_vectors_idx = (t * (X_petal.dot(w) + b) < 1) # find the locations
115 # of the miss classified data points that fall within the vectors
116 sv = X_petal[support_vectors_idx]
117
118 plt.scatter(sv[:, 0], sv[:, 1], s=180,marker='o', facecolors='none',edgecolors='k')
119
120 %% compute the confusion matrix and F1 score
121
122 y_predicted = scaled_svm_clf.predict(X_petal)
123 confusion_matrix = sk.metrics.confusion_matrix(y_predicted, y_petal)
124 f1_score = sk.metrics.f1_score(y_predicted, y_petal)
125
126 print(f1_score)

```

Example 6.2

```

1  """
2  Example 6.2 Polynomial Features
3  @author: Austin R.J. Downey
4  """
5
6  import IPython as IP
7  IP.get_ipython().run_line_magic('reset', '-sf')
8
9  import numpy as np
10 import matplotlib.pyplot as plt
11 import sklearn as sk
12 from sklearn import datasets
13 from sklearn import pipeline
14 from sklearn import svm
15
16 plt.close('all')
17
18 %% Build and plot the data
19
20 # build the data
21 X, y = sk.datasets.make_moons(n_samples=100, noise=0.25, random_state=2)
22
23 plt.figure()
24 plt.plot(X[:,0][y==0],X[:,1][y==0], 's')
25 plt.plot(X[:,0][y==1],X[:,1][y==1], 'd')
26 plt.xlabel("$x_1$")
27 plt.ylabel("$x_2$")
28
29 %% SVM polynomial features
30 svm_clf = sk.pipeline.Pipeline([
31     ("poly_features", sk.preprocessing.PolynomialFeatures(degree=3)),
32     ("scaler", sk.preprocessing.StandardScaler()),
33     ("svm_clf", sk.svm.LinearSVC(C=10))
34 ])
35 svm_clf.fit(X, y)
36
37 # make the 2d space for the color
38 x1 = np.linspace(-2, 3, 200)
39 x2 = np.linspace(-2, 2, 100)
40 x1_grid, x2_grid = np.meshgrid(x1, x2)
41
42 # calculate the binary decions and predection values
43 X2 = np.vstack((x1_grid.ravel(), x2_grid.ravel())).T
44 y_pred = svm_clf.predict(X2).reshape(x1_grid.shape)
45 y_decision = svm_clf.decision_function(X2).reshape(x1_grid.shape)
46
47 con_lines = [-30,-20,-10,-5,-2,-1,0,1,2,5,10,20,30]
48
49
50 # plot the figure
51 plt.figure()
52 # provide the solid background color for classification
53 plt.contourf(x1_grid, x2_grid, y_pred, cmap=plt.cm.brg, alpha=0.2)
54 # add the contour colors for the threshold
55 plt.contourf(x1_grid, x2_grid, y_decision, con_lines, cmap=plt.cm.brg, alpha=0.1)
56 # add the contour lines
57 contour = plt.contour(x1_grid, x2_grid, y_decision, con_lines, cmap=plt.cm.brg)
58 plt.clabel(contour, inline=1, fontsize=12)
59 plt.plot(X[:, 0][y==0], X[:, 1][y==0], "s")
60 plt.plot(X[:, 0][y==1], X[:, 1][y==1], "d")
61 plt.xlabel("$x_1$")
62 plt.ylabel("$x_2$")

```

Example 6.3

```

1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3  """
4  Example 6.3 SVM polynomial kernel
5
6  @author: Austin R.J. Downey
7  """
8
9  import IPython as IP
10 IP.get_ipython().run_line_magic('reset', '-sf')
11
12 import numpy as np
13 import matplotlib.pyplot as plt
14 import sklearn as sk
15 from sklearn import datasets
16 from sklearn import pipeline
17 from sklearn import svm
18
19 plt.close('all')
20
21
22 """ Build and plot the data
23
24 X, y = sk.datasets.make_moons(n_samples=100, noise=0.25, random_state=2)
25
26 plt.figure()
27 plt.plot(X[:,0][y==0],X[:,1][y==0], 's')
28 plt.plot(X[:,0][y==1],X[:,1][y==1], 'd')
29 plt.xlabel("$x_1$")
30 plt.ylabel("$x_2$")
31 plt.grid(True)
32
33
34 """ SVM polynomial features
35 svm_clf = sk.pipeline.Pipeline([
36     ("scaler", sk.preprocessing.StandardScaler()),
37     ("svm_clf", sk.svm.SVC(kernel="poly", degree=3, coef0=1, C=100))
38 ])
39 svm_clf.fit(X, y)
40
41
42 """ Prepare for plotting
43 # make the 2d space for the color
44 x1 = np.linspace(-2, 3, 200)
45 x2 = np.linspace(-2, 2, 100)
46 x1_grid, x2_grid = np.meshgrid(x1, x2)
47
48 # calculate the binary decisions and prediction values
49 X2 = np.vstack((x1_grid.ravel(), x2_grid.ravel())).T
50 y_decision = svm_clf.decision_function(X2).reshape(x1_grid.shape)
51
52 # plot the figure (the show up in the figure environment defined above
53 con_lines = [-30,-20,-10,-5,0,5,10,20,30]
54 contour = plt.contour(x1_grid, x2_grid, y_decision, con_lines, cmap=plt.cm.brg)
55 plt.clabel(contour, inline=1, fontsize=12)
56
57

```

Example 6.4

```

1  """
2  Example 6.4 SVM Regression
3  @author: Austin R.J. Downey
4  """
5
6  import IPython as IP
7  IP.get_ipython().run_line_magic('reset', '-sf')
8
9  import numpy as np
10 import matplotlib.pyplot as plt
11 import sklearn as sk
12 from sklearn import svm
13
14
15 plt.close('all')
16
17 """ build the data sets
18 np.random.seed(2) # 2 and 6 are pretty good
19 m = 100
20 X = 6 * np.random.rand(m,1) - 3
21 y = 0.5 * X**2 + X + 2 + np.random.randn(m,1)
22 y = y.ravel()
23
24 # plot the data
25 plt.figure()
26 plt.grid(True)
27 plt.plot(X,y,'o')
28 plt.xlabel('x')
29 plt.ylabel('y')
30
31
32 """ SVM regression
33
34 svm_reg = sk.svm.SVR(kernel="rbf", degree=3, C=1, epsilon=0.8, gamma="scale")
35 # Try poly kernel, and different degree, C, and epsilon values
36 svm_reg.fit(X, y)
37 x1 = np.linspace(-3, 3, 100).reshape(100, 1)
38 y_pred = svm_reg.predict(x1)
39
40
41 # plot the SVR model on top of the existing data
42 plt.plot(x1, y_pred, "-", linewidth=2, label=r"$\hat{y}$")
43 plt.plot(x1, y_pred + svm_reg.epsilon, "g--", label='curb')
44 plt.plot(x1, y_pred - svm_reg.epsilon, "g--")
45 plt.scatter(X[svm_reg.support_], y[svm_reg.support_], s=100, marker='o', facecolor='none',
46            edgecolors='gray')
47 plt.legend(loc="upper left")

```